

Manuscript version: Published Version

The version presented in WRAP is the published version (Version of Record).

Persistent WRAP URL:

<http://wrap.warwick.ac.uk/157205>

How to cite:

The repository item page linked to above, will contain details on accessing citation guidance from the publisher.

Copyright and reuse:

The Warwick Research Archive Portal (WRAP) makes this work by researchers of the University of Warwick available open access under the following conditions.

Copyright © and all moral rights to the version of the paper presented here belong to the individual author(s) and/or other copyright owners. To the extent reasonable and practicable the material made available in WRAP has been checked for eligibility before being made available.

Copies of full items can be used for personal research or study, educational, or not-for-profit purposes without prior permission or charge. Provided that the authors, title and full bibliographic details are credited, a hyperlink and/or URL is given for the original metadata page and the content is not changed in any way.

Publisher's statement:

Please refer to the repository item page, publisher's statement section, for further information.

For more information, please contact the WRAP Team at: wrap@warwick.ac.uk

The linear conditional expectation in Hilbert space

ILJA KLEBANOV^{1,*}, BJÖRN SPRUNGK² and T.J. SULLIVAN^{1,3,†}

¹Zuse Institute Berlin, Takustraße 7, 14195 Berlin, Germany. E-mail: ^{*}klebanov@zib.de; [†]sullivan@zib.de

²Technische Universität Bergakademie Freiberg, 09596 Freiberg, Germany.

E-mail: bjoern.sprungk@math.tu-freiberg.de

³Mathematics Institute and School of Engineering, The University of Warwick, Coventry, CV4 7AL, United Kingdom. E-mail: t.j.sullivan@warwick.ac.uk

The *linear conditional expectation* (LCE) provides a best linear (or rather, affine) estimate of the conditional expectation and hence plays an important rôle in approximate Bayesian inference, especially the *Bayes linear* approach. This article establishes the analytical properties of the LCE in an infinite-dimensional Hilbert space context. In addition, working in the space of affine Hilbert–Schmidt operators, we establish a regularisation procedure for this LCE. As an important application, we obtain a simple alternative derivation and intuitive justification of the *conditional mean embedding* formula, a concept widely used in machine learning to perform the conditioning of random variables by embedding them into reproducing kernel Hilbert spaces.

Keywords: Bayes linear analysis; conditional mean embedding; reproducing kernel Hilbert space; linear conditional expectation

1. Introduction

The crucial step in most inference problems is the approximation of the conditional expectation $\mathbb{E}[U|V]$, where $U \in L^2(\Omega, \Sigma, \mathbb{P}; \mathcal{G})$ and $V \in L^2(\Omega, \Sigma, \mathbb{P}; \mathcal{H})$ are random variables over some probability space $(\Omega, \Sigma, \mathbb{P})$ taking values in some separable Hilbert spaces \mathcal{G} and \mathcal{H} , respectively. In Bayesian statistics, where it relates to the posterior mean, $\mathbb{E}[U|V]$ is an important point estimator of the inferred parameter. It is well known¹ that $\mathbb{E}[U|V]$ is the best approximation of U by a $\sigma(V)$ -measurable random variable within $L^2(\Omega, \sigma(V), \mathbb{P}; \mathcal{G})$ (i.e. the orthogonal projection of U onto $L^2(\Omega, \sigma(V), \mathbb{P}; \mathcal{G})$),

$$\mathbb{E}[U|V] = \arg \min_{\tilde{U} \in L^2(\Omega, \sigma(V); \mathcal{G})} \|U - \tilde{U}\|_{L^2(\Omega, \Sigma, \mathbb{P}; \mathcal{G})} = \arg \min_{\tilde{U} \in L^2(\Omega, \sigma(V); \mathcal{G})} \mathbb{E}[\|U - \tilde{U}\|_{\mathcal{G}}^2]. \quad (1.1)$$

By the Doob–Dynkin representation [21], Lemma 1.13, the conditional expectation can therefore be rewritten in the form

$$\mathbb{E}[U|V] = \gamma_{U|V} \circ V \quad \mathbb{P}\text{-almost surely}, \quad (1.2)$$

where $\gamma_{U|V}: \mathcal{H} \rightarrow \mathcal{G}$ is a measurable map which we will call the *conditional expectation function* (CEF). In the language of statistical learning theory (or statistical decision theory), $\gamma_{U|V}$ is called the *regression function* and constitutes a Bayes predictor for the least squares error loss, that is, the predictor with the minimal risk [20], Section 2.4, which follows directly from (1.1).

¹For \mathbb{R} -valued random variables see, for example, [10], Theorem 10.2.9; the general case follows by choosing orthonormal bases.

While computing $\gamma_{U|V}$, which is the main object of interest, is infeasible in most applications, various estimates can be constructed. The most prominent approach is to approximate $\gamma_{U|V}$ within the class $\mathbf{A}(\mathcal{H}; \mathcal{G})$ of bounded affine operators² from \mathcal{H} to \mathcal{G} , since this provides an explicit formula for the *linear conditional expectation function* (LCEF) $\gamma_{U|V}^{\mathbf{A}}$ under appropriate conditions [12], Lemma 4.1:

$$\gamma_{U|V}^{\mathbf{A}}(v) = \mu_U + C_{UV}C_V^{-1}(v - \mu_V), \quad (1.3)$$

where μ_U and μ_V denote the means and C_{UV} and C_V denote the cross-covariance and covariance operators of U and V , as defined in Section 3.

While the *linear conditional expectation* (LCE) $\mathbb{E}^{\mathbf{A}}[U|V] := \gamma_{U|V}^{\mathbf{A}} \circ V$ (also known as *Bayes linear estimator* or *adjusted expectation*) has been discussed extensively by Michael Goldstein and his collaborators in the framework of Bayes linear statistics mostly from an application point of view [18], a rigorous mathematical analysis of the LCE is yet to be established, especially for the case of infinite-dimensional \mathcal{G} and \mathcal{H} . This level of generality, which this article seeks to provide, yields not just a satisfying mathematical theory but is also necessary for the application of LCE-type methods to problems with high-dimensional unknowns or data, such as time series and functional data analysis.

The first contribution of this paper is to fill this theoretical gap by studying the properties of the LCE and generalising formula (1.3) to infinite-dimensional Hilbert spaces. Thus far, (1.3) has been derived under the assumptions that \mathcal{H} is finite-dimensional and that C_V is invertible [12], Section 4. In addition, working in the spaces of (affine) Hilbert–Schmidt operators, we establish a rigorous justification for the regularised version of (1.3).

Our second contribution is a simple alternative derivation and intuitive explanation of the widely used formula for the conditional mean embedding (CME), a method used in machine learning to perform the conditioning of random variables by embedding them into RKHSs, where it reduces to an affine transformation similar to (1.3) [16,22,33]. This result follows almost directly from the fact that, by the reproducing property, $\mathbb{E}[U|V]$ coincides with its best affine approximation $\mathbb{E}^{\mathbf{A}}[U|V]$.

Note that this paper considers only centered (cross-)covariance operators defined by (3.1). Some, but not all, of the results can be proven similarly for uncentered operators defined by (3.2), the theory for which is less general, since it allows only for strictly linear instead of affine approximations, that is, one would be restricted to fitting the probability density or data points in Figure 1 with a straight line through the origin.

This paper is structured as follows. Section 2 briefly surveys related work in statistics, machine learning, and dynamical systems. Section 3 establishes notation and standing assumptions for the remainder of the paper. Section 4 forms the core of the paper, in which we study the rigorous generalisation of the LCE to the infinite-dimensional Hilbert space context and also consider multiple formulations of the linear conditional covariance operator. We analyse their basic properties (Theorems 4.5 and 4.7) and derive explicit formulae for them in several regimes (Theorems 4.8, 4.13, and 4.14). In Sections 5 and 6, these ideas are applied to kernel conditional mean embeddings of random variables into RKHSs and to the conditioning of infinite-dimensional Gaussian random vectors, respectively. Some closing remarks are given in Section 7, after which all the proofs of results in the main text are given in Section 8. The Appendix contains technical supporting results.

²We note here an unfortunate but seemingly unavoidable clash of terminology: while the approximate conditional expectation (1.3) is usually called the *linear* conditional expectation in the literature, it in fact corresponds to approximation using *affine* operators.

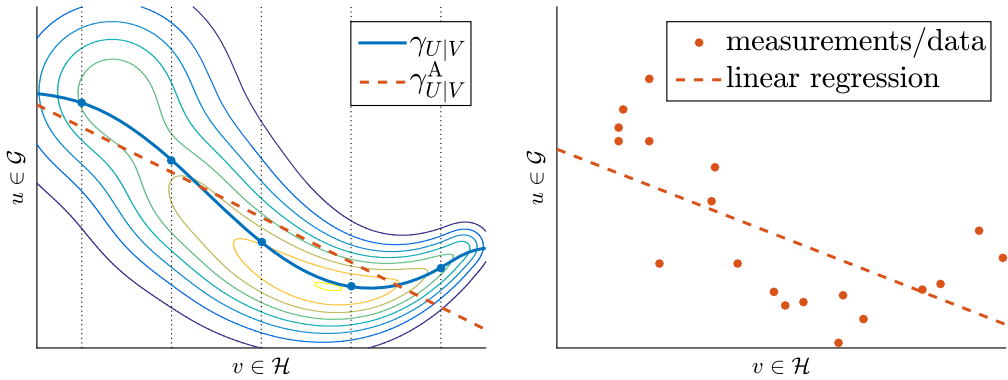


Figure 1. *Left:* Comparison of the conditional expectation function (CEF) $\gamma_{U|V} : \mathcal{H} \rightarrow \mathcal{G}$ and the linear conditional expectation function (LCEF) $\gamma_{U|V}^A \in \mathcal{A}(\mathcal{H}; \mathcal{G})$. The contour plot shows the probability density $p_{V,U}$ of (V, U) . *Right:* For an empirical probability distribution (e.g., given by data), the LCEF coincides with the solution to the linear least squares regression problem.

2. Related work

Formula (1.3) is the fundamental solution in linear least squares regression (or general linear models) and can be interpreted as the *best linear unbiased estimate* (BLUE); see, for example, [20], Sections 2.3.1 and 3.2. Figure 1 illustrates the connection between $\gamma_{U|V}^A$ and linear regression: the two coincide if the probability distribution $\mathbb{P}_{V,U}$ of (V, U) is an empirical distribution $\mathbb{P}_{V,U} = J^{-1} \sum_{j=1}^J \delta_{(v_j, u_j)}$, where (v_j, u_j) , $j = 1, \dots, J$, are (or can be thought of as) measurements or data points.

Apart from the connection to linear regression, this work is related to several fields of applied mathematics. First and foremost, it should be seen as a systematic and rigorous treatment as well as an extension of *Bayes linear analysis*, which has been introduced and investigated by Michael Goldstein and his collaborators, see, for example, [17] and [18] and the references therein. Furthermore, [34], page 9, offers a Bayesian interpretation of the BLUE in special cases, namely that in the “uninformative” infinite-variance limit of a Gaussian prior, the limiting posterior is Gaussian with the BLUE as its conditional expectation.

The LCE is applied in a variety of fields:

In geostatistics, the LCE appears in form of the Kriging estimate for the value of a random field at unexplored locations given available (noisy) data of the random field at measurement locations [6,34].

In data assimilation, formula (1.3) defines the update scheme of the Kálmán filter and its many variants, including the ensemble Kálmán filter [13]. Although this update rule is typically interpreted as a Gaussian approximation – “in the large ensemble size limit the EnKF [...] does not reproduce the filtering distribution, except in the linear Gaussian case” [31] – it has been argued by [12], Section 4, that it should rather be seen as the best linear approximation of the required conditional expectations.

In machine learning, the method of *conditional mean embedding* (CME; [16,33]) applies the conditioning formula (1.3) to random variables embedded into RKHSs, where it becomes exact (i.e., $\gamma_{U|V}^A = \gamma_{U|V}$) under certain conditions; see [22]. Section 5 provides an alternative derivation of the CME formula based on linear conditional expectations and thereby a natural justification of CMEs based in BLUEs; to the best of our knowledge, this connection has not been made before.

In the field of dynamical systems, the LCE is an important estimate of the *Koopman operator*

$$\mathcal{K}_\tau : L^\infty(\mathcal{X}) \rightarrow L^\infty(\mathcal{X}), \quad \mathcal{K}_\tau f(x) := \mathbb{E}[f(X_{t+\tau}) | X_t = x],$$

of a time-homogeneous Markov chain $(X_t)_{t \in \mathcal{T}}$, where $\tau > 0$, $\mathcal{X} \subseteq \mathbb{R}^d$, and \mathcal{T} is typically \mathbb{Z} , \mathbb{N} , \mathbb{R} , or $[0, \infty)$. Independently of one another, the dynamical systems, fluid dynamics, and molecular dynamics communities have developed data-driven dimensionality reduction techniques based on the eigendecomposition of this approximation, resulting in the methods of *time-lagged independent component analysis* (TICA) and *dynamic mode decomposition* (DMD). More generally, the data can be transformed to some feature space $\tilde{\mathcal{X}}$ by a feature map $\psi: \mathcal{X} \rightarrow \tilde{\mathcal{X}}$ prior to computing the linear approximation, resulting in the *variational approach of conformation dynamics* (VAC) or *empirical dynamic mode decomposition* (EDMD). The connections among these methods have been discussed in detail by [25]; see also the references therein to the original papers on these methods. As in the case of CMEs, if the feature map ψ is the canonical feature map $\psi(x) = k(x, \cdot)$ of an RKHS \mathcal{H} with reproducing kernel k , then our analysis is applicable and it reveals the exactness of conditioning formula (1.3), thereby annihilating one of the main sources of error – the other being the approximation of the means and covariance operators μ_U , μ_V , C_V , and C_{UV} . The resulting kernelised versions of VAC and EDMD have been studied by [32] and [26]; the exactness of (1.3), which we establish, strengthens the analytical power of these methods. We also mention that there are time-inhomogeneous variants of these methods, namely *coherent mode decomposition* (CMD) and the *variational approach for Markov processes* (VAMP) and their kernelised version *kernel canonical correlation analysis* (kernel CCA; [24]); again, the established exactness of formula (1.3) in RKHSs can be exploited for analysing these approaches.

3. Preliminaries and notation

This paper will make use of various standing assumptions and items of notation, which we collect here for easy reference.

Throughout, $(\mathcal{F}, \langle \cdot, \cdot \rangle_{\mathcal{F}})$, $(\mathcal{G}, \langle \cdot, \cdot \rangle_{\mathcal{G}})$, and $(\mathcal{H}, \langle \cdot, \cdot \rangle_{\mathcal{H}})$ will be separable real Hilbert spaces and $U \in L^2(\Omega, \Sigma, \mathbb{P}; \mathcal{G})$, $V \in L^2(\Omega, \Sigma, \mathbb{P}; \mathcal{H})$, and $W \in L^2(\Omega, \Sigma, \mathbb{P}; \mathcal{F})$, where $(\Omega, \Sigma, \mathbb{P})$ is a fixed probability space.

\mathcal{G} -valued expected values $\mu_U := \mathbb{E}[U] := \int_{\Omega} U(\omega) d\mathbb{P}(\omega)$ are always meant in the sense of a Bochner integral [8], Chapter II, as are the cross-covariance operators

$$C_{UV} := \text{Cov}[U, V] := \mathbb{E}[(U - \mathbb{E}[U]) \otimes (V - \mathbb{E}[V])] = \mathbb{E}[U \otimes V] - \mathbb{E}[U] \otimes \mathbb{E}[V] \quad (3.1)$$

from \mathcal{H} into \mathcal{G} , where, for $h \in \mathcal{H}$ and $g \in \mathcal{G}$, the outer product $g \otimes h: \mathcal{H} \rightarrow \mathcal{G}$ is the rank-one linear operator $(g \otimes h)(h') := \langle h, h' \rangle_{\mathcal{H}} g$. Naturally, we write $C_U = \text{Cov}[U]$ for the covariance operator $\text{Cov}[U, U]$, which is self-adjoint, non-negative, and trace-class [3,30], and all of the above reduces to the usual definitions in the scalar-valued case. Using Theorem A.3 and [28], Lemmas 16.7 and 16.21, it follows that the cross-covariance operators C_{UV} and $C_{VU} = C_{UV}^*$ are also trace-class (and in particular Hilbert–Schmidt) operators. In Section 5, we will briefly consider *uncentred* (cross-)covariance operators

$${}^u C_{UV} := {}^u \text{Cov}[U, V] := \mathbb{E}[U \otimes V], \quad {}^u C_U := {}^u \text{Cov}[U] := {}^u \text{Cov}[U, U]. \quad (3.2)$$

The orthogonal projection onto a closed linear subspace F of a Hilbert space \mathcal{H} will be denoted by $P_F^{\mathcal{H}}$, or just P_F whenever \mathcal{H} is clear from context. Further, we abbreviate $L^2(\Omega, \Sigma, \mathbb{P}; \mathcal{G})$ by $L^2(\mathbb{P}; \mathcal{G})$ and further by $L^2(\mathbb{P})$ if $\mathcal{G} = \mathbb{R}$; and $L^2(\Omega, \tilde{\Sigma}, \mathbb{P}_{\tilde{\Sigma}}; \mathcal{G})$ by $L^2(\Omega, \tilde{\Sigma}; \mathcal{G})$ for any sub- σ -algebra $\tilde{\Sigma} \subseteq \Sigma$. \mathbb{P}_X denotes the distribution of a random variable $X: \Omega \rightarrow \mathcal{X}$, that is, the pushforward $X_{\#}\mathbb{P}$ of \mathbb{P} under X .

For a linear operator $A: \mathcal{H} \rightarrow \mathcal{G}$ between Hilbert spaces \mathcal{H} and \mathcal{G} , its Moore–Penrose pseudo-inverse $A^\dagger: \text{dom } A^\dagger \rightarrow \mathcal{H}$ is the unique extension of

$$A|_{(\ker A)^\perp}^{-1}: \text{ran } A \rightarrow (\ker A)^\perp$$

to a linear operator A^\dagger defined on $\text{dom } A^\dagger := (\text{ran } A) \oplus (\text{ran } A)^\perp \subseteq \mathcal{G}$ subject to the criterion that $\ker A^\dagger = (\text{ran } A)^\perp$. In general, $\text{dom } A^\dagger$ is a dense but proper subspace of \mathcal{G} and A^\dagger is an unbounded operator; global definition and boundedness of A^\dagger occur precisely when $\text{ran } A$ is closed in \mathcal{G} [11], Section 2.1.

The following spaces of linear and affine operators from \mathcal{H} to \mathcal{G} will play a fundamental rôle in the approximation of $\gamma_U|_V$.

Definition 3.1. Let \mathcal{H} , \mathcal{G} , and V be as above. We define the following spaces of linear and affine operators from \mathcal{H} to \mathcal{G} :

$$\begin{aligned} \mathcal{L}(\mathcal{H}; \mathcal{G}) &:= \{\gamma: \mathcal{H} \rightarrow \mathcal{G} \mid \gamma \text{ is a bounded linear operator}\}, \\ \mathcal{A}(\mathcal{H}; \mathcal{G}) &:= \{\gamma: \mathcal{H} \rightarrow \mathcal{G} \mid \gamma(h) = b + Ah \text{ for some } b \in \mathcal{G}, A \in \mathcal{L}(\mathcal{H}; \mathcal{G})\}, \\ \mathcal{L}_V(\mathcal{H}; \mathcal{G}) &:= \{\gamma: \mathcal{H} \rightarrow \mathcal{G} \mid \gamma \text{ is linear and } \gamma \circ V \in L^2(\mathbb{P}; \mathcal{G})\}, \\ \mathcal{A}_V(\mathcal{H}; \mathcal{G}) &:= \{\gamma: \mathcal{H} \rightarrow \mathcal{G} \mid \gamma(h) = b + Ah \text{ for some } b \in \mathcal{G}, A \in \mathcal{L}_V(\mathcal{H}; \mathcal{G})\}, \\ \mathcal{L}_2(\mathcal{H}; \mathcal{G}) &:= \{\gamma: \mathcal{H} \rightarrow \mathcal{G} \mid \gamma \text{ is a Hilbert–Schmidt operator}\}, \\ \mathcal{A}_2(\mathcal{H}; \mathcal{G}) &:= \{\gamma: \mathcal{H} \rightarrow \mathcal{G} \mid \gamma(h) = b + Ah \text{ for some } b \in \mathcal{G}, A \in \mathcal{L}_2(\mathcal{H}; \mathcal{G})\}. \end{aligned}$$

Note well that elements of $\mathcal{L}_V(\mathcal{H}; \mathcal{G})$ and $\mathcal{A}_V(\mathcal{H}; \mathcal{G})$ may be unbounded operators, although their unboundedness is in some sense restricted by the square-integrability requirement. For any collection Γ of affine or linear operators $\gamma: \mathcal{H} \rightarrow \mathcal{G}$, we set $\Gamma \circ V := \{\gamma \circ V \mid \gamma \in \Gamma\}$ and

$$\overline{\mathcal{L}_G \circ V} := \overline{\mathcal{L}(\mathcal{H}; \mathcal{G}) \circ V}^{L^2(\Omega, \Sigma, \mathbb{P}; \mathcal{G})}, \quad \overline{\mathcal{A}_G \circ V} := \overline{\mathcal{A}(\mathcal{H}; \mathcal{G}) \circ V}^{L^2(\Omega, \Sigma, \mathbb{P}; \mathcal{G})}.$$

Here and henceforth, overlines and superscripts denote topological closures. The operator norm will be denoted by $\|\cdot\|$. For any affine operator $\gamma \in \mathcal{A}_V(\mathcal{H}; \mathcal{G})$, $\gamma(h) = b + Ah$, $b \in \mathcal{G}$, $A \in \mathcal{L}_V(\mathcal{H}; \mathcal{G})$, we define the “non-affine part” by $\overline{\gamma} := A$. The Hilbert–Schmidt inner product will be denoted by $\langle \gamma_1, \gamma_2 \rangle_{\mathcal{L}_2} := \text{tr}(\gamma_1^* \gamma_2) = \text{tr}(\gamma_1 \gamma_2^*)$, where $\gamma_1, \gamma_2 \in \mathcal{L}_2(\mathcal{H}; \mathcal{G})$, and the corresponding norm by $\|\cdot\|_{\mathcal{L}_2}$. Further, for $\gamma, \gamma' \in \mathcal{A}_2(\mathcal{H}; \mathcal{G})$, we define the seminorm $\|\gamma\|_{\mathcal{A}_2} := \|\overline{\gamma}\|_{\mathcal{L}_2}$ and the semi-inner product $\langle \gamma, \gamma' \rangle_{\mathcal{A}_2} := \langle \overline{\gamma}, \overline{\gamma'} \rangle_{\mathcal{L}_2}$.

Proposition 3.2. With the notation above,

$$\overline{\mathcal{L}_G \circ V} \subseteq \mathcal{L}_V(\mathcal{H}; \mathcal{G}) \circ V, \quad \overline{\mathcal{A}_G \circ V} \subseteq \mathcal{A}_V(\mathcal{H}; \mathcal{G}) \circ V.$$

4. Linear conditional expectation and covariance

It is well known that the conditional expectation $\mathbb{E}[U|V]$ is the orthogonal projection of U onto $L^2(\Omega, \sigma(V), \mathbb{P}; \mathcal{G})$; see Footnote 1. Since $\mathbb{E}[U|V]$ is $\sigma(V)$ -measurable, the Doob–Dynkin lemma

[21], Lemma 1.13, implies the existence of a Borel-measurable function $\gamma_{U|V} : \mathcal{H} \rightarrow \mathcal{G}$ such that $\mathbb{E}[U|V] = \gamma_{U|V} \circ V$ a.s. In particular, $\gamma_{U|V}$ minimizes the functional

$$\mathcal{E}_{U|V}(\gamma) := \|U - \gamma \circ V\|_{L^2(\Omega, \Sigma, \mathbb{P}; \mathcal{G})}^2 = \mathbb{E}[\|U - \gamma \circ V\|_{\mathcal{G}}^2] \quad (4.1)$$

within the class of Borel-measurable functions $\gamma : \mathcal{H} \rightarrow \mathcal{G}$. Since $\mathbb{E}[U|V]$ is unique (as an orthogonal projection), $\gamma_{U|V}$ is unique \mathbb{P}_V -a.e. and we set $\mathbb{E}[U|V = v] := \gamma_{U|V}(v)$ for $v \in \mathcal{H}$.

It seems natural to define the best linear approximation (see Footnote 2) of the conditional expectation as $\mathbb{E}^A[U|V] = \gamma_{U|V}^A \circ V$, where $\gamma_{U|V}^A$ minimizes $\mathcal{E}_{U|V}(\gamma)$ within the class $\mathbf{A}(\mathcal{H}; \mathcal{G})$ of bounded affine operators, in other words, as the $L^2(\mathbb{P}; \mathcal{G})$ -orthogonal projection of U onto $\mathbf{A}(\mathcal{H}; \mathcal{G}) \circ V$. Since this space is not closed in $L^2(\mathbb{P}; \mathcal{G})$, the proper definition uses the projection onto its closure. In line with the definition of the conditional covariance operator

$$\text{Cov}[U, W|V] := \mathbb{E}[R[U|V] \otimes R[W|V] | V], \quad R[U|V] := U - \mathbb{E}[U|V], \quad (4.2)$$

we further define the linear conditional covariance operator $\text{Cov}^A[U, W|V]$ as follows.

Definition 4.1. With the notation of Section 3, define the *linear conditional expectation* (LCE) $\mathbb{E}^A[U|V]$ (also called *adjusted expectation*, [18], Section 3.1), the *linear conditional residual* $R^A[U|V]$, and the *linear conditional covariance operator* (LCC) $\text{Cov}^A[U, W|V]$ of U given V by

$$\begin{aligned} \mathbb{E}^A[U|V] &:= P_{\mathbf{A}_G \circ V} U, \\ R^A[U|V] &:= U - \mathbb{E}^A[U|V], \\ \text{Cov}^A[U, W|V] &:= \mathbb{E}^A[R^A[U|V] \otimes R^A[W|V] | V]. \end{aligned}$$

Further, define the *average linear conditional covariance operator* (ALCC) by

$$\text{Cov}_V^A[U, W] := \mathbb{E}[R^A[U|V] \otimes R^A[W|V]].$$

By Proposition 3.2, $\mathbb{E}^A[U|V]$ will be of the form $\gamma_{U|V}^A \circ V$, where $\gamma_{U|V}^A \in \mathbf{A}_V(\mathcal{H}; \mathcal{G})$ is unique \mathbb{P}_V -a.e. and will be referred to as the *linear conditional expectation function* (LCEF). As usual, we set $\text{Cov}[U|V] := \text{Cov}[U, U|V]$, $\text{Cov}^A[U|V] := \text{Cov}^A[U, U|V]$, and $\text{Cov}_V^A[U] := \text{Cov}_V^A[U, U]$.

Since $\mathbb{E}^A[U|V]$ and $\text{Cov}^A[U|V]$ are defined as $L^2(\mathbb{P}; \mathcal{G})$ -orthogonal projection, all statements and identities in the following subsections only hold \mathbb{P} -a.s.

Remark 4.2. [18], Section 3.3, call $\text{Cov}_V^A[U, W]$ the *adjusted covariance* and argue that this is the proper way to define the linear analogue of the conditional covariance. However, this definition is not in line with the classical conditional covariance (4.2) because it fails to condition on V a second time; see Example 4.3 below. Note that the ALCC $\text{Cov}_V^A[U, W]$ is therefore not a random variable but rather the expected value of the LCC $\text{Cov}^A[U, W|V]$, hence our term “average linear conditional covariance”; see Theorem 4.15, where we also show that it coincides with the well-known Gaussian conditional covariance formula, $\text{Cov}_V^A[U, W] = C_{UW} - C_{UV} C_V^\dagger C_{VW}$ (and similarly its more general version in the incompatible case). While $\text{Cov}_V^A[U, U]$ is always non-negative (see Theorem 4.15(b)), $\text{Cov}^A[U, W|V]$ can take on negative values (see Theorem 4.7(i)).

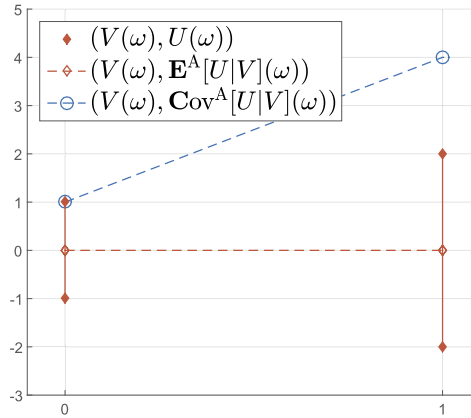


Figure 2. In Example 4.3, the conditional expectation $\mathbb{E}[U|V]$ as well as the conditional covariance $\text{Cov}[U|V]$ happen to be affine, which is why they coincide with the LCE $\mathbb{E}^A[U|V]$ and the LCC $\text{Cov}^A[U|V]$, respectively. The ALCC $\text{Cov}_V^A[U] = \frac{5}{2}$ captures the expected value of $\text{Cov}^A[U|V]$.

Example 4.3. Consider the following simple example of an LCE and an LCC. Let $\mathcal{H} = \mathcal{G} = \mathbb{R}$, let \mathbb{P} be the uniform distribution on $\Omega = \{1, 2, 3, 4\}$, and let V and U be as defined below and illustrated in Figure 2.

$\omega \in \Omega$	$V(\omega)$	$U(\omega)$	$\mathbb{E}^A[U V](\omega)$	$\text{Cov}^A[U V](\omega)$
1	0	1	0	1
2	0	-1	0	1
3	1	2	0	4
4	1	-2	0	4

By symmetry, $\mathbb{E}^A[U|V] = \mathbb{E}[U|V] = 0$ and, since $(V, R^A[U|V]^2)$ takes on only two values, $\text{Cov}^A[U|V]$ coincides with the classical conditional covariance $\text{Cov}[U|V]$. The ALCC $\text{Cov}_V^A[U] = \frac{5}{2}$ only captures its expected value.

The main aim of this section is to investigate basic properties of, and provide explicit formulae for, the LCE and the LCC.

4.1. Basic properties of the LCE

This section highlights, in Theorems 4.5 and 4.7, respectively, the ways in which the LCE shares and lacks the key properties of the exact conditional expectation. We call attention to the non-trivial conditions that appear to be necessary for the LCE version of the dominated convergence theorem; see Theorem 4.5(i), Remark 4.6, and Theorem 4.7(g). As mentioned above, all statements concerning $\mathbb{E}^A[U|V]$ and $\text{Cov}^A[U|V]$ only hold \mathbb{P} -a.s.

Lemma 4.4. *With the notation of Section 3, let $W \in \overline{\mathbf{A}_{\mathcal{F}} \circ V}$. Then, a.s.,*

$$\mathbb{E}[R^A[U|V]] = 0, \quad \text{Cov}[R^A[U|V], W] = 0.$$

In particular,

$$\text{Cov}[\mathbb{E}^A[U|V], V] = \text{Cov}[U, V] \quad \text{a.s.}$$

By way of comparison with the exact conditional expectation, some basic properties of the LCE are summarised by the following theorem (a martingale property together with a martingale convergence theorem are postponed to Theorem 4.12).

Theorem 4.5 (Basic properties satisfied by the LCE and the LCC). *With the notation of Section 3, let U' and $U_k \in L^2(\Omega, \Sigma, \mathbb{P}; \mathcal{G})$ for $k \in \mathbb{N}$. Further, let $\varphi \in \mathcal{A}(\mathcal{H}; \mathcal{F})$. The LCE fulfills the following basic properties \mathbb{P} -a.s.:*

- (a) stability:
 $\mathbb{E}^A[U|V] = \mathbb{E}[U|V]$, if $\mathbb{E}[U|V] \in \overline{\mathcal{A}_G \circ V}$, in particular,
 $\mathbb{E}^A[U|V] = g$ a.s. whenever $U = g \in \mathcal{G}$ a.s., $\mathbb{E}^A[V|V] = V$ and $\mathbb{E}^A[\varphi \circ V|V] = \varphi \circ V$;
- (b) linearity:
 $\mathbb{E}^A[aU + bU'|V] = a\mathbb{E}^A[U|V] + b\mathbb{E}^A[U'|V]$ for any $a, b \in \mathbb{R}$ and
 $\mathbb{E}^A[\psi(U)|V] = \psi(\mathbb{E}^A[U|V])$ for any $\psi \in \mathcal{A}(\mathcal{G}; \mathcal{F})$;
- (c) self-adjointness:
 $\mathbb{E}[\langle U', \mathbb{E}^A[U|V] \rangle_{\mathcal{G}}] = \mathbb{E}[\langle \mathbb{E}^A[U'|V], \mathbb{E}^A[U|V] \rangle_{\mathcal{G}}] = \mathbb{E}[\langle \mathbb{E}^A[U'|V], U \rangle_{\mathcal{G}}]$;
- (d) law of total linear expectation:
 $\mathbb{E}[\mathbb{E}^A[U|V]] = \mathbb{E}[U]$;
- (e) compatibility with conditional expectation:
 $\mathbb{E}[\mathbb{E}^A[U|V] | W] = \mathbb{E}^A[\mathbb{E}[U|W] | V]$;
 $\mathbb{E}[\mathbb{E}^A[U|V] | V] = \mathbb{E}^A[\mathbb{E}[U|V] | V] = \mathbb{E}^A[U|V]$;
- (f) tower properties:
 $\mathbb{E}^A[\mathbb{E}[U|V] | W] = \mathbb{E}^A[U|W]$ if $\sigma(W) \subseteq \sigma(V)$;
 $\mathbb{E}^A[U | \mathbb{E}^A[U|V]] = \mathbb{E}^A[U|V]$;
 $\mathbb{E}^A[\mathbb{E}^A[U|V] | \varphi \circ V] = \mathbb{E}^A[\mathbb{E}[U|V] | \varphi \circ V] = \mathbb{E}^A[U | \varphi \circ V]$, in particular,
 $\mathbb{E}^A[\mathbb{E}^A[U|(V, W)] | V] = \mathbb{E}^A[\mathbb{E}[U|(V, W)] | V] = \mathbb{E}^A[U|V]$;
- (g) law of total linear covariance:
 $\text{Cov}[U, W] = \text{Cov}[\mathbb{E}^A[U|V], \mathbb{E}^A[W|V]] + \mathbb{E}[\text{Cov}^A[U, W|V]]$, in particular,
 $\text{Cov}[U] \geq \text{Cov}[\mathbb{E}[U|V]] \geq \text{Cov}[\mathbb{E}^A[U|V]] \geq 0$.
- (h) pulling out independent factors:
 $\mathbb{E}^A[W \otimes U|V] = \mathbb{E}[W] \otimes \mathbb{E}^A[U|V]$, if W is independent of (U, V) , in particular,
 $\mathbb{E}^A[W|V] = \mathbb{E}[W]$, if V and W are independent (this also follows from (a));
- (i) L^2 -dominated convergence theorem (DCT):
 $\|\mathbb{E}^A[U_k|V] - \mathbb{E}^A[U|V]\|_{\mathcal{G}} \xrightarrow[k \rightarrow \infty]{\text{a.s.}} 0$ if C_V has finite rank and either of the following conditions holds:
 $(\alpha) \|U_k - U\|_{\mathcal{G}} \xrightarrow[k \rightarrow \infty]{\text{a.s.}} 0$ and $\|U_k\|_{\mathcal{G}} \leq Y$ for all $k \in \mathbb{N}$ and some $Y \in L^2(\mathbb{P})$,
 $(\beta) \|U_k - U\|_{L^2(\mathbb{P}; \mathcal{G})} \xrightarrow[k \rightarrow \infty]{} 0$.

Remark 4.6 (Sufficient conditions for the DCT). Note that in Theorem 4.5(i)(α) the dominating random variable $Y \in L^2(\mathbb{P})$ is assumed to be square-integrable, which is slightly stronger than the conventional assumption $Y \in L^1(\mathbb{P})$ (a counterexample to the sufficiency of the latter assumption is provided in Theorem 4.7(g)). On the other hand, (β) is a particularly weak assumption, too weak for analogous statements on the (regular) conditional expectation $\mathbb{E}[\cdot|V]$ in place of $\mathbb{E}^A[\cdot|V]$.

So far, we could only prove the DCT under the condition that C_V has finite rank. Note that counterexamples can easily be constructed (see below) if one only assumes (β) . However, the validity of the DCT under (α) for C_V of infinite rank remains an open problem. One obstacle here is the missing monotonicity of the LCE, see Theorem 4.7(a): $\|U_k\|_{\mathcal{G}} \leq Y$ a.s. does not imply that $\|\mathbb{E}^A[U_k|V]\|_{\mathcal{G}} \leq Y$ a.s.

For a counterexample to the DCT under assumption (β) (without the finite-rank assumption) consider a centered Gaussian random variable V on $\mathcal{H} = \ell^2$ with Karhunen–Loève expansion $V = \sum_n \sigma_n Z_n e_n$, where $\sigma_n > 0$ for all $n \in \mathbb{N}$, $\sum_n \sigma_n^2 < \infty$, $Z_n \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1)$ and $(e_n)_{n \in \mathbb{N}}$ is the canonical basis of $\mathcal{H} = \ell^2$. Choose $\mathcal{G} = \mathbb{R}$ and $U_k = \delta_k Z_k$ where $\delta_k \searrow 0$ such that $\mathbb{P}[A_k] = 1/k$ for $A_k = [\delta_k Z_k \geq \varepsilon]$ and $\varepsilon = 1$. Then $\mathbb{E}^A[U_k|V] = U_k$ by definition of the LCE and, since the family of events $(A_k)_{k \in \mathbb{N}}$ is independent and $\sum_k \mathbb{P}[A_k] = \infty$, the second Borel–Cantelli lemma implies $\mathbb{P}[A_k \text{ i.o.}] = 1$. Hence, by [4], Theorem 4.1.6, $\mathbb{E}^A[U_k|V] = U_k$ does not converge to zero a.s., while (β) is satisfied since $\delta_k \rightarrow 0$.

It is also worth mentioning that the LCE does *not* fulfill several important properties of conditional expectations. These are summarised by the following statements and the actual counterexamples are provided in the proof in Section 8. Note in particular that there are scalar-valued counterexamples in each case, and so these deficiencies of the LCE are not merely a consequence of the Hilbert space context.

Theorem 4.7 (Basic properties *not* satisfied by the LCE and the LCC). *For each of the following desired properties of the LCE, there is an explicit counterexample using \mathbb{R} -valued random variables U , U' , V , etc. satisfying the assumptions of Theorem 4.5. As usual, all statements have to be understood \mathbb{P} -a.s.*

- (a) monotonicity:^(invalid)
 $U \geq U'$ implies $\mathbb{E}^A[U|V] \geq \mathbb{E}^A[U'|V]$ in the case $\mathcal{G} = \mathbb{R}$;
- (b) triangle inequality:^(invalid)
 $\|\mathbb{E}^A[U|V]\|_{\mathcal{G}} \leq \mathbb{E}^A[\|U\|_{\mathcal{G}}|V]$;
- (c) Jensen's inequality:^(invalid)
 $f(\mathbb{E}^A[U|V]) \leq \mathbb{E}^A[f(U)|V]$ for any convex³ function $f: \mathcal{G} \rightarrow \mathbb{R}$;
- (d) pulling out known factors:^(invalid)
 $\mathbb{E}^A[f(V)U|V] = f(V)\mathbb{E}^A[U|V]$ for measurable⁴ maps $f: \mathcal{H} \rightarrow \mathbb{R}$;
- (e) (yet another) tower property:^(invalid)
 $\mathbb{E}^A[\mathbb{E}^A[U|V] | W] = \mathbb{E}^A[U|W]$ if $\sigma(W) \subseteq \sigma(V)$;
- (f) Fatou's lemma:^(invalid)
 Let $\mathcal{G} = \mathbb{R}$ and $\mathbb{E}^A[\inf_{k \in \mathbb{N}} U_k|V] < \infty$ (alternatively, $\mathbb{E}[\inf_{k \in \mathbb{N}} U_k|V] < \infty$). Then
 $\mathbb{E}^A[\liminf_{k \rightarrow \infty} U_k | V] \leq \liminf_{k \rightarrow \infty} \mathbb{E}^A[U_k|V]$;
- (g) L^1 -dominated convergence theorem:^(invalid)
 If $\|U_k - U\|_{\mathcal{G}} \xrightarrow[k \rightarrow \infty]{\text{a.s.}} 0$ and $\|U_k\|_{\mathcal{G}} \leq Y$ for all $k \in \mathbb{N}$ and some $Y \in L^1(\mathbb{P})$, then $\|\mathbb{E}^A[U_k|V] - \mathbb{E}^A[U|V]\|_{\mathcal{G}} \xrightarrow[k \rightarrow \infty]{\text{a.s.}} 0$.
- (h) L^p -contractivity for $p \neq 2$:^(invalid)
 $\mathbb{E}^A[\cdot|V]$ is a contractive projection of $L^p(\mathbb{P}; \mathcal{G})$ spaces for some $1 \leq p \neq 2$, i.e.
 $\mathbb{E}[\|\mathbb{E}^A[U|V]\|_{\mathcal{G}}^p] \leq \mathbb{E}[\|U\|_{\mathcal{G}}^p]$ for any $U \in L^p(\Omega, \Sigma, \mathbb{P}; \mathcal{G})$.

³Note that Jensen's (in-)equality holds for affine functions $f \in \mathcal{A}(\mathcal{G}; \mathcal{F})$, cf. Theorem 4.5(b).

⁴Our counterexample shows that property (d) is invalid even if “measurable” is strengthened to “linear”.

(i) non-negativity of the LCC:^(invalid)

The LCC $\mathbb{C}\text{ov}^A[U|V]$ is non-negative, $\mathbb{C}\text{ov}^A[U|V] \geq 0$.

Note that, in contrast to Theorem 4.7(h), $\mathbb{E}^A[\cdot|V]$ is a contractive projection on $L^2(\mathbb{P}; \mathcal{G})$; this follows directly from the definition of the LCE as an $L^2(\mathbb{P}; \mathcal{G})$ -orthogonal projection.

4.2. Explicit formula for the LCE: Compatible case

We are first going to assume that $\text{ran } C_{VU} \subseteq \text{ran } C_V$, which, following [7] and [29], we call the *compatible case*. In this case, the orthogonal projection $\mathbb{E}^A[U|V]$ of U onto $A_{\mathcal{G}} \circ V$ turns out to lie in $A(\mathcal{H}; \mathcal{G}) \circ V$, which is generally not closed in $L^2(\Omega, \Sigma, \mathbb{P}; \mathcal{G})$. The following theorem provides an explicit formula for the (affine) conditional mean and generalises [12], Lemma 4.1.

Theorem 4.8 (Formula for the LCE: compatible case). *With the notation of Section 3, if the range inclusion $\text{ran } C_{VU} \subseteq \text{ran } C_V$ holds, then the operator $C_V^\dagger C_{VU} : \mathcal{G} \rightarrow \mathcal{H}$ is bounded and the operator $\gamma_{U|V}^A \in A(\mathcal{H}; \mathcal{G})$ defined by*

$$\gamma_{U|V}^A(v) := \mu_U + (C_V^\dagger C_{VU})^*(v - \mu_V)$$

minimizes the functional $\mathcal{E}_{U|V}$ given by (4.1) within $A(\mathcal{H}; \mathcal{G})$. In particular, $\mathbb{E}^A[U|V] = \gamma_{U|V}^A \circ V$ a.s., that is, $\gamma_{U|V}^A$ is an LCEF.

Theorem 4.8 is a genuine generalisation of the case $\dim \mathcal{H} < \infty$ for the following reason, which is a direct consequence of Theorem A.3.

Corollary 4.9. *With the notation of Section 3, the condition $\text{ran } C_{VU} \subseteq \text{ran } C_V$ is always fulfilled whenever C_V has closed range, and, in particular, if \mathcal{H} is finite dimensional.*

4.3. Explicit formula for the LCE: Incompatible case

We are now going to treat the general case, in which the orthogonal projection $\mathbb{E}^A[U|V]$ of U can not be expected to lie in $A(\mathcal{H}; \mathcal{G}) \circ V$. We are therefore going to approximate $\mathbb{E}^A[U|V]$ by a sequence of bounded (in fact, even finite-rank) operators $\gamma_{U|V}^{(n)} \in A(\mathcal{H}; \mathcal{G})$ composed with V , where the convergence will hold in the $L^2(\mathbb{P}; \mathcal{G})$ norm as well as a.s. This requires some additional notation.

Notation 4.10. Let $\dim \mathcal{H} = \infty^5$ and recall the notation of Section 3. Consider the eigendecomposition of the covariance operator C_V ,

$$C_V = \sum_{n \in \mathbb{N}} \sigma_n h_n \otimes h_n, \quad \sigma_n^2 \geq 0,$$

⁵This assumption is not substantial and we make it merely for the sake of simplifying our notation. Note that the finite-dimensional case has been analysed in the previous subsection.

where $(h_n)_{n \in \mathbb{N}}$ is a complete orthonormal system of \mathcal{H} . Let $n \in \mathbb{N}$, $\mathcal{H}^{(n)} := \text{span}\{h_1, \dots, h_n\}$, $V^{(n)} := P_{\mathcal{H}^{(n)}} V$ and

$$C := \begin{pmatrix} C_U & C_{UV} \\ C_{VU} & C_V \end{pmatrix}, \quad C^{(n)} := P_{\mathcal{H}^{(n)}} C P_{\mathcal{H}^{(n)}} = \begin{pmatrix} C_U & C_{UV}^{(n)} \\ C_{VU}^{(n)} & C_V^{(n)} \end{pmatrix}.$$

Since $C_V^{(n)}$ has finite rank, Theorem A.3 yields $\text{ran } C_{VU}^{(n)} \subseteq \text{ran } C_V^{(n)}$ and we can define the operator $\gamma_{U|V}^{(n)} \in \mathcal{A}(\mathcal{H}; \mathcal{G})$ by

$$\gamma_{U|V}^{(n)}(v) := \mu_U + (C_V^{(n)})^\dagger C_{VU}^{(n)}(v - \mu_V).$$

Further, by Theorem A.3 and adopting the notation therein, the operator $M_{VU} := (C_V^{1/2})^\dagger C_{VU} = R_{VU} C_U^{1/2}: \mathcal{G} \rightarrow \mathcal{H}$ is well defined and bounded (in fact, it is even Hilbert–Schmidt).

Lemma 4.11. *With the notation of Section 3 and using Notation 4.10,*

$$\overline{A_{\mathcal{G}} \circ V} \cap L^2(\Omega, \sigma(V^{(n)}); \mathcal{G}) = \overline{A_{\mathcal{G}} \circ V^{(n)}}. \quad (4.3)$$

Theorem 4.12 (Martingale property and martingale convergence theorem). *With the notation of Section 3 and using Notation 4.10,*

(a) $(\mathbb{E}^A[U|V^{(n)}])_{n \in \mathbb{N}}$ is a martingale with respect to the filtration $(\sigma(V^{(n)}))_{n \in \mathbb{N}}$; more precisely,

$$\mathbb{E}[\mathbb{E}^A[U|V] | V^{(n)}] = \mathbb{E}^A[U|V^{(n)}] \quad \text{a.s.}; \quad (4.4)$$

(b) the following martingale convergence theorem holds a.s. and in $L^p(\mathbb{P}; \mathcal{G})$ for $1 \leq p \leq 2$:

$$\mathbb{E}^A[U|V^{(n)}] \xrightarrow{n \rightarrow \infty} \mathbb{E}^A[U|V]. \quad (4.5)$$

Theorem 4.13 (Formula for the LCE: incompatible case). *With the notation of Section 3 and using Notation 4.10, the operators $\gamma_{U|V}^{(n)}$ minimize the functional $\mathcal{E}_{U|V}$ given by (4.1) within $\mathcal{A}(\mathcal{H}; \mathcal{G})$ for $n \rightarrow \infty$, i.e.*

$$\mathcal{E}_{U|V}(\gamma_{U|V}^{(n)}) \xrightarrow{n \rightarrow \infty} \inf_{\gamma \in \mathcal{A}(\mathcal{H}; \mathcal{G})} \mathcal{E}_{U|V}(\gamma).$$

In other words,

$$\|\mathbb{E}^A[U|V] - \gamma_{U|V}^{(n)} \circ V\|_{L^2(\Omega, \Sigma, \mathbb{P}; \mathcal{G})} \xrightarrow{n \rightarrow \infty} 0. \quad (4.6)$$

Further, denoting the pushforward of \mathbb{P} under V by \mathbb{P}_V , $\gamma_{U|V}^{(n)} \circ V$ converges to $\mathbb{E}^A[U|V]$ a.s.,

$$\|\gamma_{U|V}^{(n)}(v) - \mathbb{E}^A[U|V = v]\|_{\mathcal{G}} \xrightarrow{n \rightarrow \infty} 0 \quad \text{for } \mathbb{P}_V\text{-a.e. } v \in \mathcal{H}. \quad (4.7)$$

4.4. Explicit formula for the LCE: Regularised case

In most practical applications the means and (cross-)covariance operators of U and V are not accessible explicitly, but have to be approximated empirically from data (in the simplest case, from independent and identically distributed samples $(u_n, v_n) \sim \mathbb{P}_{UV}$, where \mathbb{P}_{UV} denotes the joint distribution of U and

V). Since the Moore–Penrose pseudo-inverse C_V^\dagger shows an unstable behaviour when approximated empirically [22], Section SM2, it is typically replaced by its regularised version $(C_V + \varepsilon \text{Id}_{\mathcal{H}})^{-1}$, where $\varepsilon > 0$ is a regularisation parameter. The following theorem shows that this is a principled way to address this issue, since the resulting operators minimize a perturbed functional $\mathcal{E}_{U|V}^{\text{reg}}$. The natural space of operators in this context turns out to be the space of affine Hilbert–Schmidt operators.

Theorem 4.14 (Formula for the LCE: regularised case). *With the notation of Section 3, the operator $\gamma_\varepsilon^{\text{A}_2} \in \text{A}_2(\mathcal{H}; \mathcal{G})$ defined by*

$$\gamma_\varepsilon^{\text{A}_2}(v) := \mu_U + C_{UV}(C_V + \varepsilon \text{Id}_{\mathcal{H}})^{-1}(v - \mu_V)$$

minimizes the Tikhonov–Philipps-regularised functional

$$\mathcal{E}_{U|V}^{\text{reg}}(\gamma) = \mathcal{E}_{U|V}(\gamma) + \varepsilon \|\gamma\|_{\text{A}_2}^2$$

within $\text{A}_2(\mathcal{H}; \mathcal{G})$, where $\mathcal{E}_{U|V}$ is given by (4.1).

4.5. Explicit formula for the linear conditional covariance

Before we derive a formula for the LCC $\text{Cov}^{\text{A}}[U, W|V]$, the following theorem states an explicit formula for and some basic properties of the ALCC $\text{Cov}_V^{\text{A}}[U, W]$. In particular, as mentioned earlier in Remark 4.2, we characterise the ALCC as the expected value $\mathbb{E}[\text{Cov}^{\text{A}}[U, W|V]]$ of the LCC – hence our terminology for each of these conditional covariances.

Theorem 4.15 (Properties of the ALCC). *With the notation of Section 3 and using Notation 4.10,*

- (a) $\mathbb{E}[\text{Cov}^{\text{A}}[U, W|V]] = \text{Cov}_V^{\text{A}}[U, W]$;
- (b) $\text{Cov}_V^{\text{A}}[U] = \mathbb{E}[\text{Cov}^{\text{A}}[U|V]] \geq \mathbb{E}[\text{Cov}[U|V]] \geq 0$;
- (c) $\text{Cov}_V^{\text{A}}[U, W] = C_{UW} - M_{VU}^* M_{VW}$,
in particular, in the compatible case $\text{ran } C_{VW} \subseteq \text{ran } C_V$,
 $\text{Cov}_V^{\text{A}}[U, W] = C_{UW} - C_{UV} C_V^\dagger C_{VW}$.

Remark 4.16. The inequality $\text{Cov}_V^{\text{A}}[U] \geq 0$ can also be seen more directly using (c):

$$\text{Cov}_V^{\text{A}}[U] = C_U - M_{VU}^* M_{VU} = C_U - (R_{VU} C_U^{1/2})^* R_{VU} C_U^{1/2} = C_U^{1/2} (\text{Id}_{\mathcal{G}} - R_{VU}^* R_{VU}) C_U^{1/2} \geq 0,$$

where we used Notation 4.10 and $\|R_{VU}\| \leq 1$ (cf. Theorem A.3). While its expected value is always non-negative, the LCC $\text{Cov}^{\text{A}}[U|V]$ itself can take on negative values, see Theorem 4.7(i).

Remark 4.17. Computing the conditional covariance $\text{Cov}[U|V = v]$ for various $v \in \mathcal{H}$ is often too costly, in which case one can focus on its mean $\mathbb{E}[\text{Cov}[U|V]]$. The above statements show that the average LCC $\mathbb{E}[\text{Cov}^{\text{A}}[U|V]] = \text{Cov}_V^{\text{A}}[U]$, which can be computed by the Gaussian conditional covariance formula, never underestimates the true expected conditional covariance $\mathbb{E}[\text{Cov}[U|V]]$.

As a consequence, we obtain an explicit formula for the LCC $\text{Cov}^{\text{A}}[U, W|V]$:

Corollary 4.18 (Formula for the LCC). *With the notation of Section 3, let*

$$Z := R^{\text{A}}[U|V] \otimes R^{\text{A}}[W|V]: \Omega \rightarrow \text{L}_2(\mathcal{F}; \mathcal{G}).$$

Then, using Notation 4.10,

$$\mathbb{C}\text{ov}^A[U, W|V] = C_{UW} - M_{VU}^* M_{VW} + \lim_{n \rightarrow \infty} (C_V^{(n)\dagger} C_{VZ}^{(n)})^* (V - \mu_V) \quad a.s.,$$

where the limit is in the Hilbert–Schmidt norm and

$$\begin{aligned} \mu_Z &= C_{UW} - M_{VU}^* M_{VW}, \\ C_{VZ} &= \mathbb{E}[V \otimes (U - \mathbb{E}^A[U|V]) \otimes (W - \mathbb{E}^A[W|V])] - \mu_V \otimes \mu_Z. \end{aligned}$$

In particular, in the compatible case with $\text{ran } C_{VW} \subseteq \text{ran } C_V$ and $\text{ran } C_{VZ} \subseteq \text{ran } C_V$,

$$\mathbb{C}\text{ov}^A[U, W|V] = C_{UW} - C_{UV} C_V^\dagger C_{VW} + (C_V^\dagger C_{VZ})^* (V - \mu_V) \quad a.s.$$

5. Application to kernel conditional mean embeddings

The above results have a beautiful application to the derivation of conditional mean embeddings (CMEs), a concept used in machine learning to perform conditioning of random variables after embedding them into suitable reproducing kernel Hilbert spaces (RKHSs). To this end, let \mathcal{H} and \mathcal{G} be two RKHSs over measurable spaces \mathcal{X} and \mathcal{Y} respectively, with reproducing kernels k and ℓ and canonical feature maps $\varphi(x) := k(x, \cdot)$ and $\psi(y) := \ell(y, \cdot)$.

For two random variables $X: \Omega \rightarrow \mathcal{X}$ and $Y: \Omega \rightarrow \mathcal{Y}$ with joint distribution \mathbb{P}_{XY} and corresponding marginal distributions \mathbb{P}_X and \mathbb{P}_Y such that $V := \varphi(X) \in L^2(\Omega, \Sigma, \mathbb{P}; \mathcal{H})$ and $U := \psi(Y) \in L^2(\Omega, \Sigma, \mathbb{P}; \mathcal{G})$, respectively, the CME $\mathbb{E}[U|X]$ can be characterised by the linear-algebraic transformation

$$\mathbb{E}[U|X] = \mu_U + (C_V^\dagger C_{VU})^* (\varphi(X) - \mu_V), \quad (5.1)$$

which holds under appropriate technical assumptions [22]. Formula (5.1) can be interpreted as saying that application of the Kálmán update or BLUE formulae to RKHS embeddings of random variables realises the embedding of conditional distributions.

The theory on (affine) linear conditional means from Section 4 provides an alternative proof and a more insightful and explanatory derivation of (5.1). The main idea is to find conditions under which the conditional expectation $\mathbb{E}[U|V]$ agrees with the linear conditional expectation $\mathbb{E}^A[U|V]$. Assuming φ to be injective implies that $\mathbb{E}[U|X] = \mathbb{E}[U|V]$ and (5.1) then follows directly from Theorem 4.8, while Theorem 4.13 implies the more generally applicable formula in Theorem 5.11. The reason why one could hope for $\mathbb{E}[U|V] = \mathbb{E}^A[U|V]$ to hold is the celebrated *kernel trick*, the guiding theme of RKHS-based methods: many nonlinear problems in the original spaces \mathcal{X} and \mathcal{Y} (here, conditioning) become linear-algebraic problems when embedded into the corresponding RKHSs \mathcal{H} and \mathcal{G} .

5.1. Setup and notation

Here, with apologies for the large notational overhead relative to the brevity of the results in Section 5.2, we introduce the precise technical assumptions and notation needed for the validity of the CME approach; see [22], Section 2, for a detailed exposition.

Regarding the kernel mean embedding of random variables $X: \Omega \rightarrow \mathcal{X}$ and $Y: \Omega \rightarrow \mathcal{Y}$ into RKHSs \mathcal{H} and \mathcal{G} over \mathcal{X} and \mathcal{Y} , respectively, with reproducing kernels k and ℓ via the canonical feature maps $\varphi(x) := k(x, \cdot)$ and $\psi(y) := \ell(y, \cdot)$ we make the following basic assumptions:

Assumption 5.1.

- (a) The space \mathcal{X} is a measurable space and \mathcal{Y} is a Borel space.
- (b) The kernel functions $k: \mathcal{X} \times \mathcal{X} \rightarrow \mathbb{R}$ and $\ell: \mathcal{Y} \times \mathcal{Y} \rightarrow \mathbb{R}$ are symmetric positive definite and measurable and the corresponding RKHSs $(\mathcal{H}, \langle \cdot, \cdot \rangle_{\mathcal{H}})$ and $(\mathcal{G}, \langle \cdot, \cdot \rangle_{\mathcal{G}})$ are separable. Moreover, the canonical feature map $\varphi(x) := k(x, \cdot)$ is injective.⁶
- (c) The random variables $V := \varphi(X)$ and $U := \psi(Y)$ lie in $L^2(\Omega, \Sigma, \mathbb{P}; \mathcal{G})$ and $L^2(\Omega, \Sigma, \mathbb{P}; \mathcal{H})$, respectively.
- (d) For any $h \in \mathcal{H}$ we have $\|h\|_{\mathcal{H}} = 0$ if and only if $h = 0$ \mathbb{P}_X -a.e. in \mathcal{X} .

By [21], Theorem 5.3, Assumption 5.1(a) ensures the existence of a \mathbb{P}_X -a.e.-unique regular version of the conditional probability distribution $\mathbb{P}_{Y|X=x}$, $x \in \mathcal{X}$, for random variables $X: \Omega \rightarrow \mathcal{X}$ and $Y: \Omega \rightarrow \mathcal{Y}$. Moreover, by [35], Lemma 4.25, Assumption 5.1(b) guarantees the measurability of $\varphi: \mathcal{X} \rightarrow \mathcal{H}$ and $\psi: \mathcal{Y} \rightarrow \mathcal{G}$, respectively. Assumption 5.1(c) implies that \mathcal{H} (resp. \mathcal{G}) is continuously embedded in the pre-Hilbert space $\mathcal{L}^2(\mathbb{P}_X)$ (resp. $\mathcal{L}^2(\mathbb{P}_Y)$). Furthermore, it also follows that $\mathbb{E}[\|\psi(Y)\|_{\mathcal{G}}^2 | X = x] < \infty$ and that \mathcal{G} is continuously embedded in $\mathcal{L}^2(\mathbb{P}_{Y|X=x})$ for all $x \in \mathcal{X}_Y$, where $\mathcal{X}_Y \subseteq \mathcal{X}$ has full \mathbb{P}_X measure, see [22], Section 2. Assumption 5.1(d) clearly holds if k is continuous and $\text{supp}(\mathbb{P}_X) = \mathcal{X}$. In particular, it allows us to view \mathcal{H} as a subspace of the Lebesgue Hilbert space $L^2(\mathbb{P}_X)$.

Subsequently, we work with Bochner spaces $L^2(\mathbb{P}_X; \mathcal{F})$ where \mathcal{F} denotes another separable Hilbert space (which in our case will be equal to either \mathbb{R} or \mathcal{G}). Recall that the space $L^2(\mathbb{P}_X; \mathcal{F})$ is isometrically isomorphic to the Hilbert tensor product space $\mathcal{F} \otimes L^2(\mathbb{P}_X)$. We comment on the various perspectives on tensor product space exploited in our proofs in Remark 5.3 below. For stating our second main result, we require the following definitions.

Notation 5.2.

- (a) Given a separable Hilbert space \mathcal{F} we define $L^2_{\mathcal{C}}(\mathbb{P}_X; \mathcal{F})$ to be the quotient space $L^2(\mathbb{P}_X; \mathcal{F})/\mathcal{C}$, where

$$\mathcal{C} := \{f \in L^2(\mathbb{P}_X; \mathcal{F}) \mid \exists c \in \mathcal{F}: f(x) = c \text{ for } \mathbb{P}_X\text{-a.e. } x \in \mathcal{X}\},$$

$$([f_1], [f_2])_{L^2_{\mathcal{C}}(\mathbb{P}_X; \mathcal{F})} := \langle f_1 - \mathbb{E}[f_1(X)], f_2 - \mathbb{E}[f_2(X)] \rangle_{L^2(\mathbb{P}_X; \mathcal{F})}.$$

In the case $\mathcal{F} = \mathbb{R}$, we abbreviate the space $L^2_{\mathcal{C}}(\mathbb{P}_X; \mathbb{R})$ by $L^2_{\mathcal{C}}(\mathbb{P}_X)$ and for any subspace $\mathcal{U} \subseteq L^2(\mathbb{P}_X; \mathcal{F})$ we define $\mathcal{U}_{\mathcal{C}} := \mathcal{U}/(\mathcal{U} \cap \mathcal{C})$ and identify it with a subspace of $L^2_{\mathcal{C}}(\mathbb{P}_X; \mathcal{F})$.

- (b) Furthermore, we define the main object of our interest, the conditional mean

$$\mathbf{m}: \mathcal{X} \rightarrow \mathcal{G}, \quad \mathbf{m}(x) := \begin{cases} \mathbb{E}[U|X=x], & \text{for } x \in \mathcal{X}_Y, \\ 0, & \text{otherwise.} \end{cases}$$

Note that $\mathbf{m} \in L^2(\mathbb{P}_X; \mathcal{G})$, since Jensen's inequality, the law of total expectation, and Assumption 5.1(c) together yield that

$$\|\mathbf{m}\|_{L^2(\mathbb{P}_X; \mathcal{G})}^2 = \mathbb{E}[\|\mathbb{E}[\psi(Y)|X]\|_{\mathcal{G}}^2] \leq \mathbb{E}[\mathbb{E}[\|U\|_{\mathcal{G}}^2|X]] = \mathbb{E}[\|U\|_{\mathcal{G}}^2] < \infty.$$

⁶The injectivity of φ was not required for the derivations in [22]. However, it represents a minor restriction since one typically considers characteristic kernels, which implies injectivity of φ .

(c) We also introduce the notation

$$f_g(x) := \mathbb{E}[g(Y)|X=x] = \langle g, \mathbf{m}(x) \rangle_{\mathcal{G}},$$

mainly for the comparison of our formulations to [22].

Remark 5.3. For $\mathcal{F} = \mathcal{H}$ and $\mathcal{F} = L^2(\mathbb{P}_X)$, we will view the Hilbert tensor product space $\mathcal{G} \otimes \mathcal{F}$ from three⁷ perspectives:

- $\mathcal{G} \otimes \mathcal{F}$ is isometrically isomorphic to the space $\mathcal{L}_2(\mathcal{F}; \mathcal{G})$ of Hilbert–Schmidt operators from \mathcal{F} to \mathcal{G} [2], Chapter 12, and we sometimes view $g \otimes f \in \mathcal{G} \otimes \mathcal{F} \cong \mathcal{L}_2(\mathcal{F}; \mathcal{G})$ as the corresponding mapping from \mathcal{F} to \mathcal{G} given by

$$[g \otimes f]_{\mathcal{F} \rightarrow \mathcal{G}}(f') := \langle f, f' \rangle_{\mathcal{F}} g.$$

- $\mathcal{G} \otimes \mathcal{F}$ can be viewed as a set of functions from \mathcal{X} to \mathcal{G} [2], Chapter 12. Thus, we can view $g \otimes f \in \mathcal{G} \otimes \mathcal{F} \subseteq L^2(\mathbb{P}_X; \mathcal{G})$ accordingly as a function on \mathcal{X} taking values in \mathcal{G} :

$$[g \otimes f]_{\mathcal{X} \rightarrow \mathcal{G}}(x) := f(x)g = [g \otimes f]_{\mathcal{H} \rightarrow \mathcal{G}}(\varphi(x)),$$

where the last equality holds for $\mathcal{F} = \mathcal{H}$ by the reproducing property, in which case we have established the important observation

$$\mathbf{f}_{\mathcal{X} \rightarrow \mathcal{G}}(x) = \mathbf{f}_{\mathcal{H} \rightarrow \mathcal{G}}(\varphi(x)), \quad \mathbf{f} \in \mathcal{G} \otimes \mathcal{H}, x \in \mathcal{X}. \quad (5.2)$$

Note that $\mathcal{G} \otimes \mathcal{F}$ is indeed a subspace of $L^2(\mathbb{P}_X; \mathcal{G})$, which is obvious in the case $\mathcal{F} = L^2(\mathbb{P}_X)$, and, in the case $\mathcal{F} = \mathcal{H}$, follows from

$$\|\mathbf{f}_{\mathcal{X} \rightarrow \mathcal{G}}\|_{L^2(\mathbb{P}_X; \mathcal{G})}^2 = \mathbb{E}[\|\mathbf{f}_{\mathcal{X} \rightarrow \mathcal{G}}(X)\|_{\mathcal{G}}^2] \stackrel{(5.2)}{\leq} \|\mathbf{f}_{\mathcal{H} \rightarrow \mathcal{G}}\|_{\mathcal{L}_2}^2 \mathbb{E}[\|\varphi(X)\|_{\mathcal{H}}^2] < \infty,$$

where we used Assumption 5.1(c) in the last step. This is hardly surprising, since $\mathcal{H} \subseteq L^2(\mathbb{P}_X)$ by Assumption 5.1(c), (d), but not trivial as we use the RKHS norm $\|\cdot\|_{\mathcal{H}}$ in the construction of $\mathcal{G} \otimes \mathcal{H}$, which may not agree with $\|\cdot\|_{L^2(\mathbb{P}_X)}$.

- Since tensor products are commutative up to isometric isomorphism, $\mathcal{G} \otimes \mathcal{F}$ is also isometrically isomorphic to $\mathcal{L}_2(\mathcal{G}; \mathcal{F})$ and we can analogously set

$$[g \otimes f]_{\mathcal{G} \rightarrow \mathcal{F}}(g') := \langle g, g' \rangle_{\mathcal{F}} f, \quad f \in \mathcal{F}, g, g' \in \mathcal{G}.$$

We will sometimes use the resulting identities for arbitrary $\mathbf{f} \in \mathcal{G} \otimes \mathcal{F}$: with $x \in \mathcal{X}$

$$\mathbf{f}_{\mathcal{G} \rightarrow \mathcal{F}}(g)(x) = \langle \mathbf{f}_{\mathcal{X} \rightarrow \mathcal{G}}(x), g \rangle_{\mathcal{G}}, \quad \langle \mathbf{f}, g \otimes f \rangle_{\mathcal{G} \otimes \mathcal{F}} = \langle \mathbf{f}_{\mathcal{G} \rightarrow \mathcal{F}}(g), f \rangle_{\mathcal{F}}. \quad (5.3)$$

However, we will drop the indices $\mathcal{F} \rightarrow \mathcal{G}$, $\mathcal{X} \rightarrow \mathcal{G}$ and $\mathcal{G} \rightarrow \mathcal{F}$ in the following, since it will always be clear which version we mean, whenever we apply $\mathbf{f} \in \mathcal{G} \otimes \mathcal{F}$ to some element of \mathcal{F} , \mathcal{X} , or \mathcal{G} , respectively.

⁷In fact, there is another viewpoint on $\mathcal{G} \otimes \mathcal{F}$, namely as a set of functions from $Y \times X$ to \mathbb{R} , where $(g \otimes f)(y, x) := g(y)f(x)$, in which case $\mathcal{G} \otimes \mathcal{F} \subseteq L^2(\mathbb{P}_Y \otimes \mathbb{P}_X)$.

The typical assumption for CMEs is that the functions f_g introduced in Notation 5.2(c) must lie in \mathcal{H} for all $g \in \mathcal{G}$. [22] discuss several weaker assumptions on f_g , which we are going to adopt in this paper. However, the main purpose of using f_g is that, for $g = \psi(y)$ with $y \in \mathcal{Y}$, $f_{\psi(y)} = \mu_{Y|X=\cdot}(y)$ and, in fact, all results in [22] rely solely on these special cases of f_g . It is therefore meaningful to restate these assumptions in terms of \mathbf{m} rather than f_g .

Assumption 5.4. Under Assumption 5.1 and using Notation 5.2, we introduce the following assumptions on the functions $\mathbf{m} \in L^2(\mathbb{P}_X; \mathcal{G})$:

- (A) $\mathbf{m} \in \mathcal{G} \otimes \mathcal{H}$;
- (B) $[\mathbf{m}] \in (\mathcal{G} \otimes \mathcal{H})_{\mathcal{C}}$;
- (C) $P_{\overline{(\mathcal{G} \otimes \mathcal{H})}_{\mathcal{C}}}^{L^2(\mathbb{P}_X; \mathcal{G})} [\mathbf{m}] \in (\mathcal{G} \otimes \mathcal{H})_{\mathcal{C}}$;
- (^uC) $P_{\overline{\mathcal{G} \otimes \mathcal{H}}}^{L^2(\mathbb{P}_X; \mathcal{G})} \mathbf{m} \in \mathcal{G} \otimes \mathcal{H}$;
- (A*) $\mathbf{m} \in \overline{\mathcal{G} \otimes \mathcal{H}}^{L^2(\mathbb{P}_X; \mathcal{G})}$;
- (B*) $[\mathbf{m}] \in \overline{(\mathcal{G} \otimes \mathcal{H})}_{\mathcal{C}}^{L^2(\mathbb{P}_X; \mathcal{G})}$.

Remark 5.5. Note that these assumptions are slightly stronger than the corresponding assumptions on f_g in [22], Section 3. The corresponding implications are formulated in Proposition 5.6 below. However, by Lemma 5.7, (B*) still follows from k being characteristic, providing a verifiable condition for the applicability of the corresponding CME formula (Theorem 5.11). Also, as before, (A*) follows from k being L^2 -universal. Indeed, if $\mathcal{G} \otimes \mathcal{H}$ is dense in $L^2(\mathbb{P}_X; \mathcal{G})$, then

$$\overline{\mathcal{G} \otimes \mathcal{H}}^{L^2(\mathbb{P}_X; \mathcal{G})} \subseteq \overline{\mathcal{G} \otimes \overline{\mathcal{H}}^{L^2(\mathbb{P}_X)}}^{\mathcal{G} \otimes L^2(\mathbb{P}_X)} = \overline{\mathcal{G} \otimes L^2(\mathbb{P}_X)}^{\mathcal{G} \otimes L^2(\mathbb{P}_X)} = L^2(\mathbb{P}_X; \mathcal{G}).$$

These and further relations among the conditions in Assumption 5.4 are summarised in Figure 3.

Proposition 5.6. *The conditions in Assumption 5.4 imply the corresponding assumptions on f_g in [22], Section 3, (here marked with a subscript “old”). More precisely, under Assumption 5.1 and using Notation 5.2,*

- (A) \implies (A_{old}) $f_g \in \mathcal{H}$ for each $g \in \mathcal{G}$;
- (B) \implies (B_{old}) $[f_g] \in \mathcal{H}_{\mathcal{C}}$ for each $g \in \mathcal{G}$;
- (C) \implies (C_{old}) $P_{\overline{\mathcal{H}_{\mathcal{C}}}}^{L^2(\mathbb{P}_X)} [f_g] \in \mathcal{H}_{\mathcal{C}}$ for each $g \in \mathcal{G}$;
- (^uC) \implies (^uC_{old}) $P_{\overline{\mathcal{H}}}^{L^2(\mathbb{P}_X)} f_g \in \mathcal{H}$ for each $g \in \mathcal{G}$;
- (A*) \implies (A*_{old}) $f_g \in \overline{\mathcal{H}}^{L^2(\mathbb{P}_X)}$ for each $g \in \mathcal{G}$;
- (B*) \implies (B*_{old}) $f_g \in \overline{(\mathcal{H})_{\mathcal{C}}}^{L^2(\mathbb{P}_X)}$ for each $g \in \mathcal{G}$.

Further,

- (C) $\implies \text{ran } C_{VU} \subseteq \text{ran } C_V$;
- (^uC) $\implies \text{ran } {}^u C_{VU} \subseteq \text{ran } {}^u C_V$.

Lemma 5.7. *Under Assumption 5.1 and using Notation 5.2, if k is a characteristic kernel, then $(\mathcal{G} \otimes \mathcal{H})_{\mathcal{C}}$ is dense in $L^2_{\mathcal{C}}(\mathbb{P}_X; \mathcal{G})$ and Assumption 5.4(B*) is satisfied.*

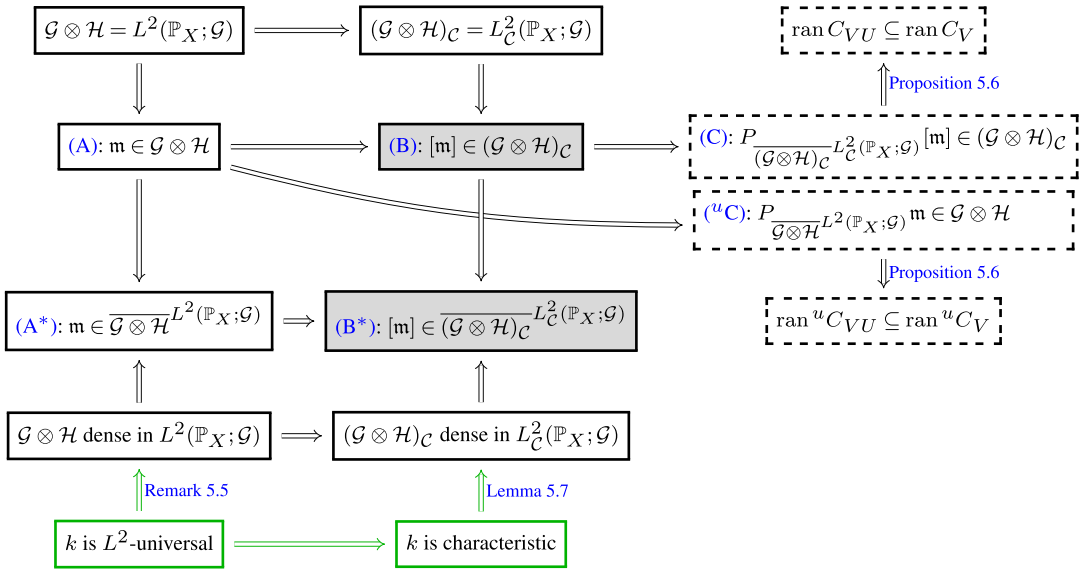


Figure 3. A hierarchy of CME-related assumptions. Sufficient conditions for validity of the CME formula are indicated by solid boxes while the insufficient Assumptions (C) and $(^u C)$, indicated by dashed boxes, have several strong theoretical implications. Assumption 5.4(B^{*}) is the most favorable one, since it is verifiable in practice, and, by Lemma 5.7, in particular is fulfilled if the kernel is universal or even just characteristic (marked in green). The shaded boxes correspond to Theorems 5.8 and 5.11.

5.2. Derivation of the CME formula

We are now in a position to rederive the CME formula under Assumption 5.4.

Theorem 5.8 (CME under Assumption 5.4(A) or (B)). Under Assumptions 5.1 and 5.4(B) the operator $C_V^\dagger C_{VU} : \mathcal{G} \rightarrow \mathcal{H}$ is bounded and

$$\mathbb{E}[U|X] = \mu_U + (C_V^\dagger C_{VU})^*(\varphi(X) - \mu_V) \quad a.s. \quad (5.4)$$

Remark 5.9. Note that we have proven a stronger statement than originally intended. Namely, the CME operator $(C_V^\dagger C_{VU})^*$ is not just bounded, but even Hilbert–Schmidt. However, it can be argued that this property is already hidden in the assumptions, namely in Assumption 5.4(B), since $\mathcal{G} \otimes \mathcal{H} \cong L_2(\mathcal{H}; \mathcal{G})$;

Remark 5.10. A similar statement can be proven for uncentered covariance operators under the stronger Assumption 5.4(A), see [22], Theorem 5.3.

Theorem 5.11 (CME under Assumption 5.4(B^{*})). Under Assumptions 5.1 and 5.4(B^{*}) and using Notation 4.10, the operators $\gamma_{U|V}^{(n)}$ satisfy

$$\begin{aligned} \|\mathbb{E}[U|X] - \gamma_{U|V}^{(n)}(\varphi(X))\|_{L^2(\mathbb{P}; \mathcal{G})} &\xrightarrow{n \rightarrow \infty} 0, \\ \|\mathbb{E}[U|X = x] - \gamma_{U|V}^{(n)}(\varphi(x))\|_{\mathcal{G}} &\xrightarrow{n \rightarrow \infty} 0 \quad \text{for } \mathbb{P}_X\text{-a.e. } x \in \mathcal{X}. \end{aligned}$$

6. Application to Gaussian conditioning in Hilbert spaces

While the conditioning of a Gaussian random variable $(U, V): \Omega \rightarrow \mathcal{G} \oplus \mathcal{H}$ on its second component is a well-established concept [19,27], the most general case (where C_V is not necessarily injective) has only been treated rather recently by [29]. In that work, by developing an approximation theory for shorted operators in terms of oblique projections and applying the martingale convergence theorem, the authors derive approximating sequences for both the conditional expectation $\mathbb{E}[U|V]$ and the conditional covariance operator $\text{Cov}[U|V]$.

The formula that [29], Theorem 3.3, obtain for the conditional expectation $\mathbb{E}[U|V]$ is identical to (4.7) (with $\mathbb{E}[U|V]$ in place of $\mathbb{E}^A[U|V]$). Similar to CMEs in Section 5, our theory provides an alternative derivation of this formula by

- (i) proving $\mathbb{E}[U|V] \in \overline{A_{\mathcal{G}} \circ \overline{V}}$, implying the identity $\mathbb{E}[U|V] = \mathbb{E}^A[U|V]$;
- (ii) applying Theorem 4.13.

Let us give a short sketch of the proof:

Proof sketch for (i). Let $\overline{U} := U - \mu_U$, $\overline{V} := V - \mu_V$ and $\gamma := \gamma_{\overline{U}|\overline{V}}: \mathcal{H} \rightarrow \mathcal{G}$ be the corresponding CEF. [36], Theorem 3.11, show that there exists a linear subspace $\tilde{\mathcal{H}}$ of \mathcal{H} such that $\overline{V} \in \tilde{\mathcal{H}}$ a.s. and the restriction $\gamma|_{\tilde{\mathcal{H}}}$ is linear. In the proof, the authors further construct a sequence $\gamma_n \in \mathcal{L}(\mathcal{H}; \mathcal{G})$ such that

$$\gamma_n(h) = \sum_{i=1}^n \langle h, h_i \rangle \gamma(h_i) \xrightarrow{n \rightarrow \infty} \gamma(h) \quad \text{for all } h \in \tilde{\mathcal{H}}.$$

Using the Karhunen–Loève expansion of \overline{V} , one can prove

$$\|\gamma_n \circ \overline{V} - \gamma \circ \overline{V}\|_{L^2(\mathbb{P}; \mathcal{G})}^2 \xrightarrow{n \rightarrow \infty} 0.$$

Hence $\mathbb{E}[\overline{U}|\overline{V}] = \gamma \circ \overline{V} \in \overline{A_{\mathcal{G}} \circ \overline{V}}$ and thereby $\mathbb{E}[U|V] = \mu_U + \gamma \circ (V - \mu_V) \in \overline{A_{\mathcal{G}} \circ \overline{V}}$. \square

The appeal to [36], Theorem 3.11, is somewhat unsatisfactory, since close inspection of the proof of that theorem reveals that it in fact establishes the entire conditional mean formula for Gaussian conditioning. In this sense, our derivation of the formula for the conditional *mean* $\mathbb{E}[U|V]$ in the Gaussian case is not novel. However, let us now turn to the conditional *covariance* $\text{Cov}[U|V]$.

It is well known that the conditional covariance is constant for Gaussian random variables (i.e. it does not depend on the value of the conditioning variable), and that, since $\mathbb{E}[U|V] = \mathbb{E}^A[U|V]$, it coincides with the ALCC, $\text{Cov}[U|V] = \text{Cov}_V^A[U]$. Therefore, our results show that, in contrast to the conditional expectation, the conditional covariance does require the approximating sequence established by [29], Theorem 3.4, but the explicit formula from Theorem 4.15(c) applies. In summary, we can consider three versions of the Gaussian conditional covariance formula:

- the *invertible case* in which C_V is invertible and, in particular, \mathcal{H} is finite dimensional:

$$\text{Cov}[U|V] = C_U - C_{UV}C_V^{-1}C_{VU}.$$

- the *compatible case* in which $\text{ran } C_{VU} \subseteq \text{ran } C_V$:

By Theorem A.1, the operator $C_V^\dagger C_{VU} \in \mathcal{L}(\mathcal{G}; \mathcal{H})$ is well defined and bounded and

$$\text{Cov}[U|V] = C_U - C_{UV}C_V^\dagger C_{VU}.$$

- the *incompatible* (or *general*) case:

By Theorem A.3, the operator $M_{VU} := (C_V^{1/2})^\dagger C_{VU}$ is well-defined and bounded and

$$\text{Cov}[U|V] = C_U - M_{VU}^* M_{VU}.$$

7. Closing remarks

This paper presents a rigorous theory of the linear conditional expectation (LCE) $\mathbb{E}^A[\cdot|V]$ that strongly extends the existing theory on Bayes linear analysis [18].

After the definitions of the linear conditional expectation $\mathbb{E}^A[U|V]$ and the linear conditional covariance (LCC) operator $\text{Cov}^A[U, W|V]$ – which is related to, but differs from, the so-called adjusted covariance used in Bayes linear statistics – we studied in detail which properties of the common conditional expectation $\mathbb{E}[\cdot|V]$ and conditional covariance $\text{Cov}[\cdot, \cdot|V]$ hold for their linear approximations. Amongst others, we proved several tower properties and the laws of total expectation and total covariance. On the other hand, $\mathbb{E}^A[\cdot|V]$ is neither monotonic, nor contractive in $L^p(\mathbb{P}; \mathcal{G})$ (except for $p = 2$, which is clear from its definition) and does not fulfill the triangle inequality. The dominated convergence theorem holds only under modified assumptions and, so far, could only be proved under the assumption that C_V has finite rank (see Theorem 4.5(i), Remark 4.6, and Theorem 4.7(g)).

We derived explicit formulae for both the LCE and the LCC, distinguishing between the so-called compatible (simple) and the incompatible (hard) case, as well as providing a regularised formula for the LCE.

Naturally, whenever $\mathbb{E}^A[U|V] = \mathbb{E}[U|V]$, these formulae apply for the common conditional expectation. This trivial observation allowed us to provide an alternative derivation of the Gaussian conditioning formulae and give a simple and intuitive proof of the widely-used technique of conditional mean embeddings (CMEs) in machine learning: it turns out that, if $U = \psi(Y)$ and $V = \varphi(X)$ are reproducing kernel Hilbert space (RKHS) embeddings of some random variables X and Y , then the above property holds true under rather mild conditions.

One direction for future work is the derivation of optimal regularisation schemes $\varepsilon(n) \rightarrow 0$ when the regularised case considered in Section 4.4 is applied to empirical sample data consisting of $n \rightarrow \infty$ data points. We anticipate that this will be a rich vein of research, but also decidedly non-trivial, since such problems admit no general solution and effective strategies (be they a priori, a posteriori, or heuristic) rely on appropriate source conditions for the unknowns.

Finally, we note that this work has concentrated on *centred* (cross-)covariance operators, which are associated with affine approximations of the conditional expectation function. Some, but not all of the statements can also be proved for *uncentred* operators, which are associated with linear approximations of the CEF and are often used in practice; for example, the uncentred formulation is commonly used for CMEs. However, the theory with uncentred operators has weaker statements and is more restrictive, and so we strongly encourage the use of centred operators.

8. Proofs

Proof of Proposition 3.2. We only give the proof of the second statement, which is similar to the first one but slightly more technical. Let $W \in \bar{A}\mathcal{G} \circ \bar{V}$ and $(\gamma_n)_{n \in \mathbb{N}}$ be a sequence in $A(\mathcal{H}; \mathcal{G})$, such that

$$\|\gamma_n \circ V - W\|_{L^2(\mathbb{P}; \mathcal{G})} \xrightarrow[n \rightarrow \infty]{\text{a.s.}} 0.$$

This implies that, for $\overline{V} := V - \mu_V$ and $\overline{W} := W - \mu_W$,

$$\|\overline{\gamma}_n \circ \overline{V} - \overline{W}\|_{L^2(\mathbb{P}; \mathcal{G})} \xrightarrow[n \rightarrow \infty]{\text{a.s.}} 0.$$

By [15], Corollary 2.32, there exists a subsequence $(\overline{\gamma}_{n_k})_{k \in \mathbb{N}}$ of $(\overline{\gamma}_n)_{n \in \mathbb{N}}$ such that

$$\|\overline{\gamma}_{n_k} \circ \overline{V}(\omega) - \overline{W}(\omega)\|_{\mathcal{G}} \xrightarrow[k \rightarrow \infty]{} 0 \quad \text{for } \mathbb{P}\text{-a.e. } \omega \in \Omega.$$

Define the linear subspace $\mathcal{H}_0 \subseteq \mathcal{H}$ and the (possibly unbounded) linear operator $A_0: \mathcal{H}_0 \rightarrow \mathcal{G}$ by

$$\mathcal{H}_0 := \{h \in \mathcal{H} \mid \overline{\gamma}_{n_k}(h) \text{ converges in } \mathcal{G}\}, \quad A_0(h) := \lim_{k \rightarrow \infty} \overline{\gamma}_{n_k}(h) \text{ for } h \in \mathcal{H}_0,$$

and extend A_0 trivially⁸ to a linear operator A on \mathcal{H} . Then $\overline{V} \in \mathcal{H}_0$ a.s. and $A \circ \overline{V} = \overline{W}$ a.s. Considering the affine operator $\gamma: \mathcal{H} \rightarrow \mathcal{G}$ given by $\gamma(h) = \mu_W + A(h - \mu_V)$ yields that $\gamma \in \mathbf{A}_V(\mathcal{H}; \mathcal{G})$ and $\gamma \circ V = W$ a.s. \square

Proof of Lemma 4.4. Since $\mathbb{E}[\mathbb{E}^A[U|V]] = \mathbb{E}[U]$ (which follows from \mathbb{E} and $\mathbb{E}^A[\cdot|V]$ being orthogonal projections, see the law of total linear expectation in Theorem 4.5(d)), it follows that $\mathbb{E}[R^A[U|V]] = 0$. Hence, by Lemma A.6, for any $\gamma \in \mathbf{L}(\mathcal{H}; \mathcal{G})$,

$$0 = \langle R^A[U|V], \gamma \circ V \rangle_{L^2(\mathbb{P}; \mathcal{G})} = \text{tr}(\text{Cov}[R^A[U|V], V] \gamma^*).$$

By Lemma A.5 this implies that $\text{Cov}[R^A[U|V], V] = 0$. Now let $W \in \overline{\mathbf{A}_{\mathcal{F}} \circ \overline{V}}$. By Proposition 3.2, $W = \gamma \circ V$ for some $\gamma \in \mathbf{A}_V(\mathcal{H}; \mathcal{F})$. Hence, invoking Lemma A.6 another time,

$$\text{Cov}[R^A[U|V], W] = (\overline{\gamma} \text{Cov}[V, R^A[U|V]])^* = 0,$$

which completes the proof. (Note that the finite trace of the cross-covariance operator was essential in the above argument.) \square

Proof of Theorem 4.5. Properties (a)–(f), except for the second statement on linearity in (b) and the second tower property in (f), follow directly from the definitions of $\mathbb{E}[U]$, $\mathbb{E}[U|V]$ and $\mathbb{E}^A[U|V]$ as orthogonal projections of U , the identity $\mathbb{E}[\langle \cdot, \cdot \rangle_{\mathcal{G}}] = \langle \cdot, \cdot \rangle_{L^2(\mathbb{P}; \mathcal{G})}$ and the inclusions

$$\mathbf{A}(\mathcal{H}; \mathcal{G}) \circ V \subseteq L^2(\sigma(V); G) \subseteq L^2(\Sigma; G), \quad \mathbf{A}(\mathcal{F}; \mathcal{G}) \circ \varphi \subseteq \mathbf{A}(\mathcal{H}; \mathcal{G}).$$

For the second statement on linearity in (b) first note that $\psi(\mathbb{E}^A[U|V]) \in \overline{\mathbf{A}_{\mathcal{F}} \circ \overline{V}}$. Lemmas A.6 and 4.4 imply that, for any $\gamma \in \mathbf{A}(\mathcal{H}; \mathcal{F})$,

$$\langle \psi(U) - \psi(\mathbb{E}^A[U|V]), \gamma \circ V \rangle_{L^2(\mathbb{P}; \mathcal{F})} = \text{tr}(\psi \text{Cov}[\mathbb{E}^A[U|V], V] \overline{\gamma}^* - \psi \text{Cov}[U, V] \overline{\gamma}^*) = 0,$$

which completes the proof of (b).

For second tower property in (f) let $\mathbb{E}^A[U|V] = \gamma \circ V \in \overline{\mathbf{A}_{\mathcal{G}} \circ \overline{V}}$ (using Proposition 3.2) and assume that there exists $\delta \in \mathbf{A}(\mathcal{G}; \mathcal{G})$ such that

$$\|U - \delta \circ \mathbb{E}^A[U|V]\|_{L^2(\mathbb{P}; \mathcal{G})} < \|U - \mathbb{E}^A[U|V]\|_{L^2(\mathbb{P}; \mathcal{G})}.$$

⁸Choose a Hamel basis \mathcal{B}_1 of \mathcal{H}_0 , extend it to a basis $\mathcal{B}_1 \cup \mathcal{B}_2$ of \mathcal{H} and set $A := \tilde{A}$ on \mathcal{B}_1 and $A := 0$ on \mathcal{B}_2 .

Then $\delta \circ \gamma \circ V \in \overline{A_{\mathcal{G}} \circ V}$ is a better $L^2(\mathbb{P}; \mathcal{G})$ -approximation of U than $\mathbb{E}^A[U|V]$, which contradicts the definition of $\mathbb{E}^A[U|V]$.

For the law of total linear covariance (g), first note that, by the law of total linear expectation (d) and Lemma 4.4,

$$\mathbb{E}[\text{Cov}^A[U, W|V]] = \mathbb{E}[\mathbb{E}^A[R^A[U|V] \otimes R^A[W|V] | V]] = \text{Cov}[R^A[U|V], R^A[W|V]].$$

Hence, again by Lemma 4.4,

$$\begin{aligned} \text{Cov}[U, W] &= \text{Cov}[\mathbb{E}^A[U|V] + R^A[U|V], \mathbb{E}^A[W|V] + R^A[W|V]] \\ &= \text{Cov}[\mathbb{E}^A[U|V], \mathbb{E}^A[W|V]] + 0 + 0 + \text{Cov}[R^A[U|V], R^A[W|V]] \\ &= \text{Cov}[\mathbb{E}^A[U|V], \mathbb{E}^A[W|V]] + \mathbb{E}[\text{Cov}^A[U, W|V]], \end{aligned}$$

proving the law of total linear covariance in (g). The inequality $\text{Cov}[U] \geq \text{Cov}[\mathbb{E}[U|V]]$ is well known (it follows from the common law of total covariance). For the second inequality, let $U' := \mathbb{E}[U|V]$. By the first tower property in (f), $\mathbb{E}^A[U'|V] = \mathbb{E}^A[U|V]$. Hence, by the law of total linear covariance that was just established,

$$\begin{aligned} \text{Cov}[U'] &= \text{Cov}[\mathbb{E}^A[U'|V]] + \mathbb{E}[\text{Cov}^A[U'|V]] = \text{Cov}[\mathbb{E}^A[U|V]] + \text{Cov}[R^A[U'|V]] \\ &\geq \text{Cov}[\mathbb{E}^A[U|V]], \end{aligned}$$

where we used Lemma 4.4 and the law of total linear expectation (d) in the second step, finalising the proof of (g).

In order to prove (h), note that the independence of W and (U, V) implies

$$\begin{aligned} \text{Cov}[W \otimes U, V] &= \mathbb{E}[W \otimes U \otimes V] - \mathbb{E}[W \otimes U] \otimes \mu_V = \mu_W \otimes \mathbb{E}[U \otimes V] - \mu_W \otimes \mu_U \otimes \mu_V \\ &= \mu_W \otimes C_{UV}. \end{aligned}$$

Further, by Lemma 4.4, $\text{Cov}[\mathbb{E}^A[U|V], V] = C_{UV}$. Since $\mathbb{E}[\mu_W \otimes \mathbb{E}^A[U|V]] = \mu_W \otimes \mu_U = \mathbb{E}[W \otimes U]$, Lemma A.6 implies that, for any $\gamma \in A(\mathcal{H}; \mathcal{F} \otimes \mathcal{G})$,

$$\begin{aligned} \langle W \otimes U - \mu_W \otimes \mathbb{E}^A[U|V], \gamma \circ V \rangle_{L^2(\mathbb{P}; \mathcal{F} \otimes \mathcal{G})} &= \text{tr}(\text{Cov}[W \otimes U - \mu_W \otimes \mathbb{E}^A[U|V], V] \bar{\gamma}^*) \\ &= \text{tr}((\text{Cov}[W \otimes U, V] - \mu_W \otimes C_{UV}) \bar{\gamma}^*) \\ &= 0. \end{aligned}$$

Therefore, $\mu_W \otimes \mathbb{E}^A[U|V]$ is the $L^2(\mathbb{P}; \mathcal{F} \otimes \mathcal{G})$ -orthogonal projection of $W \otimes U$ onto $\overline{A_{\mathcal{F} \otimes \mathcal{G}} \circ V}$.

In order to prove⁹ (i) first note that, by linearity (b), we may assume that $U = 0$. Further, since $(\alpha) \implies (\beta)$, we may simply assume that (β) holds. This implies that

$$\mu_{U_k} \xrightarrow[k \rightarrow \infty]{} 0, \quad \|C_{U_k}\| \xrightarrow[k \rightarrow \infty]{} 0, \quad C_{U_k V} = C_{U_k}^{1/2} R_{U_k V} C_V^{1/2} \xrightarrow[k \rightarrow \infty]{} 0, \quad (8.1)$$

⁹A simpler proof can be obtained from using Theorem 4.8: After establishing (8.1), (i) follows from

$$\mathbb{E}^A[U_k|V] = \mu_{U_k} + (C_V^\dagger C_{V U_k})^* (V - \mu_V) \xrightarrow[k \rightarrow \infty]{\text{a.s.}} 0.$$

where we used Theorem A.3 and adopted the notation therein (note that $C_{U_k V}$ has finite rank, since C_V has finite rank by assumption). By Lemma 4.4 and Lemma A.6,

$$\overline{\gamma}_{U_k|V}^A C_V = \text{Cov}[\gamma_{U_k|V}^A \circ V, V] = \text{Cov}[\mathbb{E}^A[U_k|V], V] = C_{U_k V} \xrightarrow[k \rightarrow \infty]{} 0.$$

Therefore, $\overline{\gamma}_{U_k|V}^A \big|_{\text{ran } C_V} \xrightarrow[k \rightarrow \infty]{} 0$ and, by the assumption of finite rank,

$$V \in \text{ran } C_V \quad \text{a.s.} \quad \text{and} \quad \overline{\gamma}_{U_k|V}^A \circ V \xrightarrow[k \rightarrow \infty]{\text{a.s.}} 0.$$

Denoting the constant part of $\gamma_{U_k|V}^A$ by $b_k \in \mathcal{G}$, that is, $\gamma_{U_k|V}^A(v) = b_k + \overline{\gamma}_{U_k|V}^A(v)$ for $v \in \mathcal{H}$, the law of total linear expectation (d) implies

$$b_k + \overline{\gamma}_{U_k|V}^A \mu_V = \mathbb{E}[\mathbb{E}^A[U_k|V]] = \mu_{U_k} \xrightarrow[k \rightarrow \infty]{\text{a.s.}} 0.$$

Since $\mu_V \in \text{ran } C_V$ and $\overline{\gamma}_{U_k|V}^A \big|_{\text{ran } C_V} \xrightarrow[k \rightarrow \infty]{} 0$, we obtain $b_k \xrightarrow[k \rightarrow \infty]{} 0$ and thereby

$$\mathbb{E}^A[U_k|V] = b_k + \overline{\gamma}_{U_k|V}^A \circ V \xrightarrow[k \rightarrow \infty]{\text{a.s.}} 0. \quad \square$$

Proof of Theorem 4.7. We choose $\mathcal{H} = \mathcal{G} = \mathcal{F} = \mathbb{R}$ for all counterexamples provided in this proof. For counterexamples to (a)–(e), let \mathbb{P} to be the uniform distribution on $\Omega = \{1, 2, 3\}$ and the random variables V , $U_1 := V$ and $U_2 := W := |V|$ be given by

$\omega \in \Omega$	$U_1(\omega) = V(\omega)$	$U_2(\omega) = W = V(\omega) $	$\mathbb{E}^A[U_1 V](\omega)$	$\mathbb{E}^A[U_2 V](\omega)$
1	−1	1	−1	$\frac{2}{3}$
2	0	0	0	$\frac{2}{3}$
3	1	1	1	$\frac{2}{3}$

Clearly, $\mathbb{E}^A[U_1|V] = U_1 = V$ and, by solving a simple linear regression (or simply by symmetry and Theorem 4.5(d)), $\mathbb{E}^A[U_2|V] \equiv \frac{2}{3}$, as illustrated in Figure 4 (left). Therefore, $U_2 \geq U_1$, but $\mathbb{E}^A[U_2|V] \not\geq \mathbb{E}^A[U_1|V]$, disproving (a). Further, $|\mathbb{E}^A[U_1|V]| \leq \mathbb{E}^A[|U_1||V]$ does not hold, providing a counterexample to (b) and (c). For $f: \mathbb{R} \rightarrow \mathbb{R}$, $f(x) = x$, we obtain $f(V)\mathbb{E}^A[U_1|V] = V^2$, which clearly cannot equal $\mathbb{E}^A[f(V)U_1|V]$, since it is not an affine transformation of V , disproving (d). Finally, $\mathbb{E}^A[\mathbb{E}^A[U_2|V] | W] = \mathbb{E}^A[\frac{2}{3}|W] = \frac{2}{3}$ a.s., which clearly differs from $\mathbb{E}^A[U_2|W] = W$ and thereby provides a counterexample to (e).

Since \mathbb{E}^A lacks monotonicity, a counterexample to (f) is easy to construct. Consider the uniform distribution \mathbb{P} on $\Omega = \{1, 2, 3\}$ and V as well as the sequence $(U_k)_{k \in \mathbb{N}}$ given by

$\omega \in \Omega$	$V(\omega)$	$U_{2k+1}(\omega)$	$U_{2k}(\omega)$	$\liminf_{k \rightarrow \infty} U_k(\omega)$	$\liminf_{k \rightarrow \infty} \mathbb{E}^A[U_k V](\omega)$
1	−1	0	1	0	$-\frac{1}{6}$
2	0	0	0	0	$\frac{2}{6}$
3	1	1	0	0	$-\frac{1}{6}$

Then $\mathbb{E}^A[\liminf_{k \rightarrow \infty} U_k|V] = \mathbb{E}^A[0|V] = 0$, while $\liminf_{k \rightarrow \infty} \mathbb{E}^A[U_k|V](\omega) < 0$ for $\omega = 1$ and $\omega = 3$, which follows from the solution of a simple linear regression problem and is visualised in Figure 4(right).

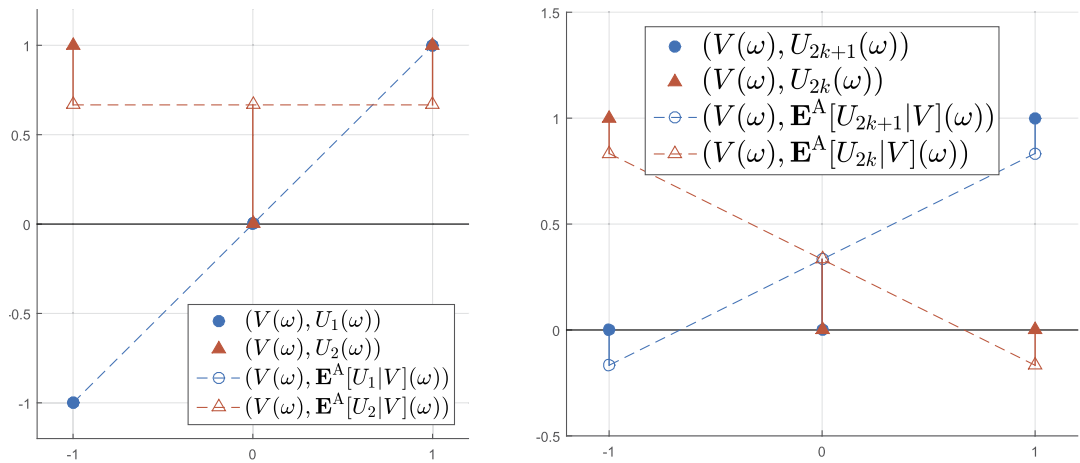


Figure 4. *Left:* Counterexample to Theorem 4.7(a)–(e) as described above, with the corresponding LCEs $\mathbb{E}^A[U_k|V]$, $k = 1, 2$. *Right:* Counterexample to Fatou's lemma (Theorem 4.7(f)) with corresponding LCEs $\mathbb{E}^A[U_k|V]$, $k \in \mathbb{N}$.

Let us now construct a counterexample to the (conventional) dominated convergence theorem (g). Let $\varepsilon > 0$ and $\alpha = (2 + 2\varepsilon)^{-1}$, for example, $\varepsilon = 1/4$ and $\alpha = 2/5$. Let \mathbb{P} be the uniform distribution on $\Omega = [-1, 1]$ and

$$V(\omega) = \begin{cases} (1 + \omega)^{-\alpha} - 1 & \text{for } \omega \in [-1, 0], \\ -(1 - \omega)^{-\alpha} + 1 & \text{for } \omega \in [0, 1], \end{cases}$$

$$U_k(\omega) = \begin{cases} (1 + \omega)^{-2\alpha} - 1 & \text{for } \omega \in \left[\frac{1}{2k} - 1, \frac{1}{k} - 1\right], \\ -(1 - \omega)^{-2\alpha} + 1 & \text{for } \omega \in \left[1 - \frac{1}{k}, 1 - \frac{1}{2k}\right], \\ 0 & \text{otherwise,} \end{cases}$$

as illustrated in Figure 5. Clearly, each U_k is bounded and thereby lies in $L^2(\mathbb{P}; \mathcal{G})$ and $U_k \xrightarrow[k \rightarrow \infty]{\text{a.s.}} 0$. Then, $\mathbb{E}^A[U_k|V] = a_k V + b_k$ for some $a_k, b_k \in \mathbb{R}$ where $b_k = 0$ for symmetry reasons. Let $\beta := 3\alpha - 1$ and note that $\beta > 0$ for sufficiently small ε ($\beta = 1/5$ in the above example). A straightforward computation shows that

$$\begin{aligned} \|a_k V - U_k\|_{L^2(\mathbb{P}; \mathcal{G})}^2 &= a_k^2 \|V\|_{L^2(\mathbb{P}; \mathcal{G})}^2 - 2a_k \langle V, U_k \rangle_{L^2(\mathbb{P}; \mathcal{G})} + \|U_k\|_{L^2(\mathbb{P}; \mathcal{G})}^2 \\ &= \frac{2(1 + \varepsilon)}{\varepsilon} a_k^2 - \frac{4k^\beta}{\beta} (2^\beta - 1) a_k + \|U_k\|_{L^2(\mathbb{P}; \mathcal{G})}^2, \end{aligned}$$

which is minimised by $a_k = \frac{(2^\beta - 1)\varepsilon}{\beta(1 + \varepsilon)} k^\beta \xrightarrow[k \rightarrow \infty]{} \infty$, which contradicts $\mathbb{E}^A[U_k|V] = a_k V \xrightarrow[k \rightarrow \infty]{\text{a.s.}} 0$.

The following two counterexamples disprove (h) for every $p < 2$ and for every $p > 2$, respectively, (for $p = 2$ the projection is clearly contractive, since it is orthogonal). Choose \mathbb{P} to be the uniform distribution on $\Omega = \{1, 2, 3, 4\}$ and the random variables V , U_1 and U_2 in the following way, where

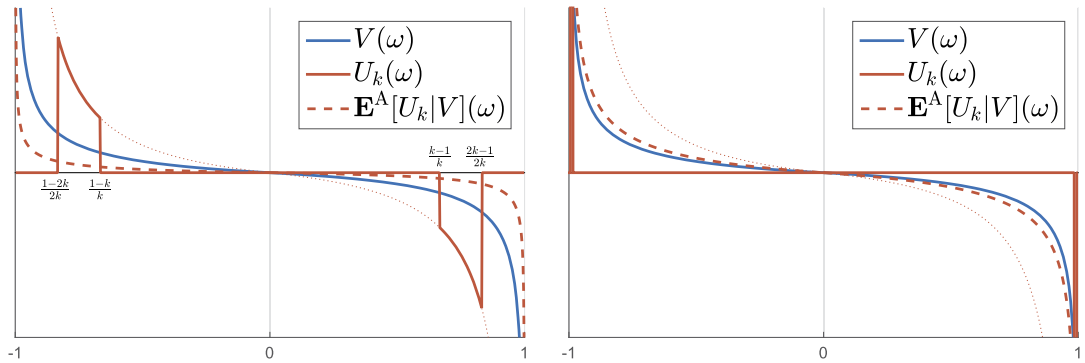


Figure 5. Counterexample to the dominated convergence theorem (Theorem 4.7(g)). Note that we plot $\mathbb{E}^A[U_k|V]$ as a function of ω (and not of V) which is why it is not a linear function in contrast to the other plots, while being a multiple of V by the factor α_k . For sufficiently small $\varepsilon > 0$, this factor α_k increases with k (here $k = 3$ (left), $k = 70$ (right) and $\varepsilon = 0.01$) as can be seen from the dotted red lines in the two plots. Therefore, in contrast to $(U_k)_{k \in \mathbb{N}}$, the sequence of LCEs $(\mathbb{E}^A[U_k|V])_{k \in \mathbb{N}}$ does not converge to zero a.s.

$\varepsilon \in (0, 1)$ is a free parameter yet to be chosen.

$\omega \in \Omega$	$V(\omega)$	$U_1(\omega)$	$U_2(\omega)$
1	-1	-1	-1
2	$-\varepsilon$	0	-2ε
3	ε	0	2ε
4	1	1	1

Again, the computation of $\mathbb{E}^A[U_j|V] = a_j V + b_j$, $a_j, b_j \in \mathbb{R}$, $j = 1, 2$, reduces to a linear regression that is solved by

$$b_1 = b_2 = 0, \quad a_1 = \frac{1}{1 + \varepsilon^2}, \quad a_2 = \frac{1 + 2\varepsilon^2}{1 + \varepsilon^2}.$$

Note that $\mathbb{E}[|U_1|^p] = \frac{1}{2}$, $\mathbb{E}[|U_2|^p] = \frac{1}{2}(1 + (2\varepsilon)^p)$ and $\mathbb{E}[\mathbb{E}^A[U_j|V]^p] = \frac{1}{2}a_j^p(1 + \varepsilon^p)$, $j = 1, 2$. It follows that the inequality in (h) for $U = U_j$, $j = 1, 2$, holds whenever

$$f_1(\varepsilon) := (1 + \varepsilon^2)^p - (1 + \varepsilon^p) \geq 0, \tag{8.2}$$

$$f_2(\varepsilon) := (1 + \varepsilon^2)^p (1 + 2^p \varepsilon^p) - (1 + 2\varepsilon^2)^p (1 + \varepsilon^p) \geq 0, \tag{8.3}$$

respectively. Bernoulli's inequality $(1 + x)^r \leq 1 + rx$ for $x \geq -1$ and exponents $0 \leq r \leq 1$ implies that, for $1 \leq p \leq 2$,

$$\begin{aligned} f_1(\varepsilon) &= (1 + \varepsilon^2)(1 + \varepsilon^2)^{p-1} - (1 + \varepsilon^p) \\ &\leq (1 + \varepsilon^2)(1 + (p-1)\varepsilon^2) - (1 + \varepsilon^p) \\ &= p\varepsilon^2 + (p-1)\varepsilon^4 - \varepsilon^p. \end{aligned}$$

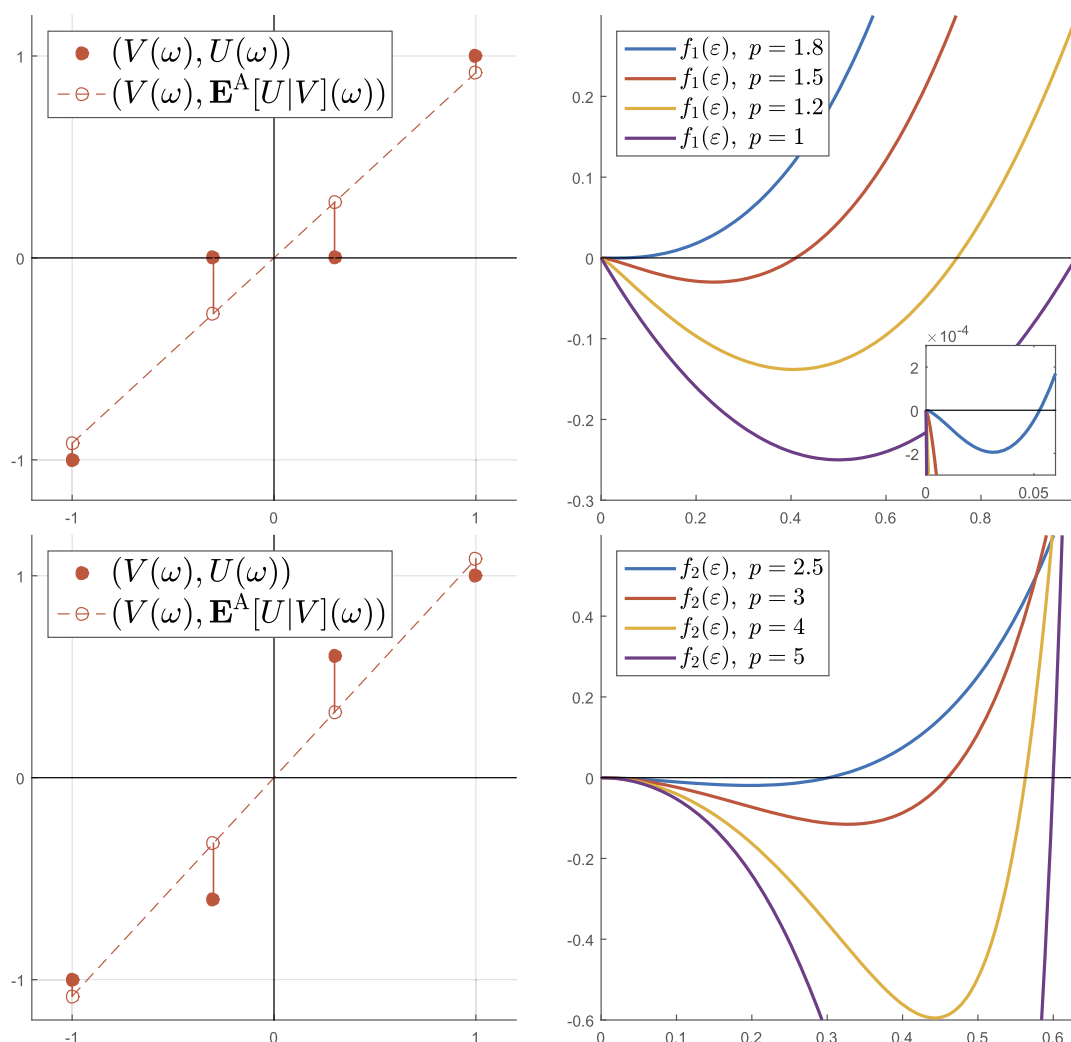


Figure 6. *Top:* Counterexample to Theorem 4.7(h) for $1 \leq p < 2$. *Bottom:* Counterexample to Theorem 4.7(h) for $p > 2$. *Left:* LCEs for the above examples and $\varepsilon = 0.3$. *Right:* The functions $f_1(\varepsilon)$ and $f_2(\varepsilon)$ for several values of p . For every $p < 2$ (respectively, $p > 2$) there exists a sufficiently small $\varepsilon > 0$ such that $f_j(\varepsilon) < 0$, $j = 1, 2$.

Since, for any $p < 2$, $p\varepsilon^2 + (p-1)\varepsilon^4 < \varepsilon^p$ for sufficiently small ε , we can falsify (8.2) and thereby disprove (h) for any $p < 2$.

For fixed $p > 2$ consider the Taylor polynomial of degree 2 for f_2 , namely $T_2 f_2(\varepsilon) = -p\varepsilon^2$; note that this is *not* a Taylor polynomial of f_2 for $p < 2$. Hence, (8.3) cannot hold for sufficiently small ε , providing a counterexample to (h) for any $p > 2$.

Figure 6 illustrates the counterexamples for $\varepsilon = 0.3$ (left) as well as the functions $f_1(\varepsilon)$ and $f_2(\varepsilon)$ for several values of p .

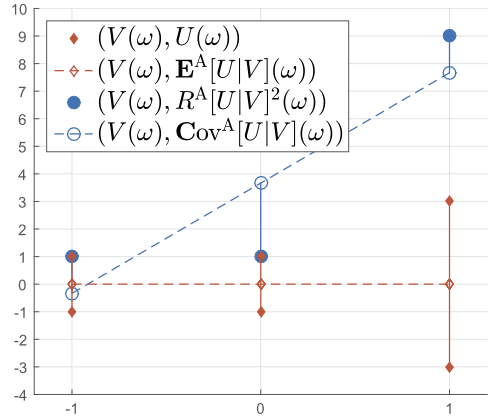


Figure 7. In contrast to the conditional covariance $\text{Cov}[U|V]$, the LCC $\text{Cov}^A[U|V]$ can take on negative values, while its expected value $\text{Cov}_V^A[U]$ is guaranteed to be non-negative.

We now give a counterexample to (i). Let \mathbb{P} be the uniform distribution on $\Omega = \{1, \dots, 6\}$ and V and U given by

$\omega \in \Omega$	$V(\omega)$	$U(\omega)$	$\mathbb{E}^A[U V](\omega)$	$R^A[U V]^2(\omega)$	$\text{Cov}^A[U V](\omega)$
1	-1	1	0	1	$\frac{1}{6}(7 - N^2)$
2	-1	-1			
3	0	1			
4	0	-1	0	1	$\frac{1}{3}(2 + N^2)$
5	1	N			
6	1	$-N$	0	N^2	$\frac{1}{6}(1 + 5N^2)$

where $N > 0$. By symmetry, $\mathbb{E}^A[U|V] = \mathbb{E}[U|V] = 0$ while $\text{Cov}^A[U|V] = \frac{1}{2}(N^2 - 1)V + \frac{1}{3}(N^2 + 2)$, which follows from the solution of a simple linear regression problem and is visualised in Figure 7. If $N > 0$ is sufficiently large, $\text{Cov}^A[U|V](\omega)$ clearly takes on negative values for $\omega = 1, 2$. \square

Proof of Theorem 4.8. First note that, by Theorem A.1, $C_V^\dagger C_{VU} \in \mathcal{L}(\mathcal{G}; \mathcal{H})$ is well-defined and bounded and that $C_V(C_V^\dagger C_{VU}) = C_{VU}$, which implies $\overline{\gamma}_{U|V}^A C_V = (C_V^\dagger C_{VU})^* C_V = C_{UV}$. We have to show that $U - \gamma_{U|V}^A \circ V$ is $L^2(\mathbb{P}; \mathcal{G})$ -perpendicular to $\gamma \circ V$ for any other $\gamma \in \mathcal{A}(\mathcal{H}; \mathcal{G})$. Since $\mathbb{E}[\gamma_{U|V}^A \circ V] = \mu_U = \mathbb{E}[U]$, it follows by Lemma A.6 that

$$\langle U - \gamma_{U|V}^A \circ V, \gamma \circ V \rangle_{L^2(\mathbb{P}; \mathcal{G})} = \text{tr}(\text{Cov}[U - \gamma_{U|V}^A \circ V, V] \overline{\gamma}^*) = \text{tr}((C_{UV} - \overline{\gamma}_{U|V}^A C_V) \overline{\gamma}^*) = 0,$$

as required. \square

Proof of Lemma 4.11. The Karhunen–Loève expansion of V takes the form

$$V = \mu_V + \sum_{i \in \mathbb{N}} Z_i h_i,$$

where the Z_i are uncorrelated real-valued random variables over $(\Omega, \Sigma, \mathbb{P})$ with $\mathbb{E}[Z_i] = 0$ and $\mathbb{V}[Z_i] = \sigma_i^2$ for all $i \in \mathbb{N}$. Observe that $\sigma(V^{(n)}) = \sigma(Z^{(n)})$, where $Z^{(n)} := (Z_1, \dots, Z_n)$.

Now let $W \in \overline{A_G \circ V} \cap L^2(\Omega, \sigma(V^{(n)}); \mathcal{G})$. Since $W \in \overline{A_G \circ V}$, there exists a sequence $(\gamma_k)_{k \in \mathbb{N}}$ in $A(\mathcal{H}; \mathcal{G})$ such that $\|\gamma_k \circ V - W\|_{L^2(\mathbb{P}; \mathcal{G})} \xrightarrow{k \rightarrow \infty} 0$. In order to show that $W \in \overline{A_G \circ V^{(n)}}$, we will find a sequence $(\tilde{\gamma}_k)_{k \in \mathbb{N}}$ in $A(\mathcal{H}; \mathcal{G})$ such that $\|\tilde{\gamma}_k \circ V^{(n)} - W\|_{L^2(\mathbb{P}; \mathcal{G})} \xrightarrow{k \rightarrow \infty} 0$. To this end, we shift each γ_k by a constant, choosing $\tilde{\gamma}_k(v) := \gamma_k(v) + \gamma_k(\mu_V - \mu_V^{(n)})$, where $\mu_V^{(n)} := \mathbb{E}[V^{(n)}] = P_{\mathcal{H}^{(n)}} \mu_V$. In order to prove above convergence observe that $\mathbb{E}[W|V^{(n)}] = W$, since W is $\sigma(V^{(n)})$ -measurable, and that

$$\mathbb{E}[\gamma_k(V - \mu_V)|V^{(n)}] = \gamma_k\left(\mathbb{E}\left[\sum_{i \in \mathbb{N}} Z_i h_i | Z^{(n)}\right]\right) = \gamma_k\left(\sum_{i=1}^n Z_i h_i\right) = \gamma_k(V^{(n)} - \mu_V^{(n)}),$$

as the random variables Z_i are uncorrelated. Since conditional expectations are $L^2(\mathbb{P}; \mathcal{G})$ -contractive projections, it follows that

$$\begin{aligned} \|\tilde{\gamma}_k \circ V^{(n)} - W\|_{L^2(\mathbb{P}; \mathcal{G})} &= \|\gamma_k(V^{(n)} - \mu_V^{(n)}) - W + \gamma_k(\mu_V)\|_{L^2(\mathbb{P}; \mathcal{G})} \\ &= \|\mathbb{E}[\gamma_k(V - \mu_V) | V^{(n)}] - \mathbb{E}[W | V^{(n)}] + \gamma_k(\mu_V)\|_{L^2(\mathbb{P}; \mathcal{G})} \\ &= \|\mathbb{E}[\gamma_k(V) - W | V^{(n)}]\|_{L^2(\mathbb{P}; \mathcal{G})} \\ &\leq \|\gamma_k \circ V - W\|_{L^2(\mathbb{P}; \mathcal{G})} \\ &\xrightarrow{k \rightarrow \infty} 0. \end{aligned}$$

The second inclusion in (4.3) is trivial. \square

Proof of Theorem 4.12. Claim (a) follows from Lemma 4.11 via

$$\mathbb{E}[\mathbb{E}^A[U|V] | V^{(n)}] = P_{L^2(\Omega, \sigma(V^{(n)}); \mathcal{G})} P_{A_G \circ V} U = P_{\overline{A_G \circ V^{(n)}}} U = \mathbb{E}^A[U|V^{(n)}],$$

where all orthogonal projections are taken with respect to the $L^2(\Omega, \Sigma, \mathbb{P}; \mathcal{G})$ inner product.

Since $\mathbb{E}^A[U|V] = \mathbb{E}[\mathbb{E}^A[U|V] | V]$ by Theorem 4.5(e) and using (4.4), claim (b) follows directly from [5], Theorems 2 and 6, or [23], Theorems 11.7 and 11.10. \square

Proof of Theorem 4.13. Since $C_V^{(n)}$ has finite rank, Theorem A.3 yields $\text{ran } C_{V|V}^{(n)} \subseteq \text{ran } C_V^{(n)}$ and Theorem 4.8 implies

$$\mathbb{E}^A[U|V^{(n)}] = P_{\overline{A_G \circ V^{(n)}}} U = \gamma_{U|V}^{(n)} \circ V^{(n)} = \gamma_{U|V}^{(n)} \circ V.$$

The statements follow directly from Theorem 4.12. \square

Proof of Theorem 4.14. First note that $\gamma_\varepsilon^{A_2} \in A_2(\mathcal{H}; \mathcal{G})$, since it is a shifted composition of the Hilbert–Schmidt operator C_{UV} and the bounded operator $(C_V + \varepsilon \text{Id}_{\mathcal{H}})^{-1}$. Now let $\gamma \in A_2(\mathcal{H}; \mathcal{G})$ and $\delta := \gamma - \gamma_\varepsilon^{A_2} \in A_2(\mathcal{H}; \mathcal{G})$. Since $\mathbb{E}[\gamma_\varepsilon^{A_2} \circ V] = \mu_U$, Lemma A.6 implies that

$$\begin{aligned} \mathbb{E}[\langle U - \gamma_\varepsilon^{A_2} \circ V, \delta \circ V \rangle_{\mathcal{G}}] &= \text{tr}(C_{UV} \bar{\delta}^* - \bar{\gamma}_\varepsilon^{A_2} C_V \bar{\delta}^*) \\ &= \text{tr}(C_{UV} (C_V + \varepsilon \text{Id}_{\mathcal{H}})^{-1} (C_V + \varepsilon \text{Id}_{\mathcal{H}} - C_V) \bar{\delta}^*) \\ &= \varepsilon \text{tr}(\bar{\gamma}_\varepsilon^{A_2} \bar{\delta}^*) \\ &= \varepsilon \langle \gamma_\varepsilon^{A_2}, \delta \rangle_{L_2}. \end{aligned}$$

Hence,

$$\begin{aligned}\mathcal{E}_{U|V}^{\text{reg}}(\gamma) &= \mathcal{E}_{U|V}^{\text{reg}}(\gamma_\varepsilon^{A_2}) + \mathbb{E}[\|\delta(V)\|_{\mathcal{G}}^2] + \varepsilon\|\delta\|_{L_2}^2 - \underbrace{2\mathbb{E}[(U - \gamma_\varepsilon^{A_2}(V), \delta(V))_{\mathcal{G}}] + 2\varepsilon(\gamma_\varepsilon^{A_2}, \delta)_{L_2}}_{=0} \\ &\geq \mathcal{E}_{U|V}^{\text{reg}}(\gamma_\varepsilon^{A_2}),\end{aligned}$$

proving the claim. \square

Proof of Theorem 4.15. By the law of total linear expectation in Theorem 4.5(d),

$$\begin{aligned}\mathbb{E}[\text{Cov}^A[U, W|V]] &= \mathbb{E}[\mathbb{E}^A[R^A[U|V] \otimes R^A[W|V] | V]] \\ &= \mathbb{E}[R^A[U|V] \otimes R^A[W|V]] \\ &= \text{Cov}_V^A[U, W],\end{aligned}$$

proving (a). By the law of total covariance and its linear version in Theorem 4.5(g), we obtain

$$\mathbb{E}[\text{Cov}[U|V]] = \text{Cov}[U] - \text{Cov}[\mathbb{E}[U|V]] \leq \text{Cov}[U] - \text{Cov}[\mathbb{E}^A[U|V]] = \mathbb{E}[\text{Cov}^A[U|V]],$$

proving (b) (the equality in (b) follows directly from (a)). In order to prove (c), first note that, by Theorem A.3 and using Notation 4.10,

$$C_V^{1/2}(\overline{\gamma}_{U|V}^{(n)})^* = C_V^{1/2}C_V^{(n)\dagger}C_{VU}^{(n)} = P_{\mathcal{H}^{(n)}}R_{VU}C_U^{1/2} \xrightarrow[n \rightarrow \infty]{} R_{VU}C_U^{1/2} = M_{VU}.$$

Hence, by (4.6) and Lemma A.6,

$$\begin{aligned}\text{Cov}[\mathbb{E}^A[U|V], \mathbb{E}^A[W|V]] &= \lim_{n \rightarrow \infty} \text{Cov}[\gamma_{U|V}^{(n)} \circ V, \gamma_{W|V}^{(n)} \circ V] \\ &= \lim_{n \rightarrow \infty} \overline{\gamma}_{U|V}^{(n)} C_V (\overline{\gamma}_{W|V}^{(n)})^* \\ &= \lim_{n \rightarrow \infty} (C_V^{1/2}(\overline{\gamma}_{U|V}^{(n)})^*)^* (C_V^{1/2}(\overline{\gamma}_{W|V}^{(n)})^*) \\ &= M_{VU}^* M_{VW}.\end{aligned}$$

By (a) and the law of total linear covariance in Theorem 4.5(g), we obtain

$$\text{Cov}_V^A[U, W] = \text{Cov}[U, W] - \text{Cov}[\mathbb{E}^A[U|V], \mathbb{E}^A[W|V]] = C_{UW} - M_{VU}^* M_{VW},$$

thus completing the proof. \square

Proof of Corollary 4.18. Noting that $\mu_Z = \text{Cov}_V^A[U, W]$, the claim follows directly from Theorems 4.8, 4.13, and 4.15. \square

Proof of Proposition 5.6. In this proof \mathfrak{m} will be viewed as an element of $L_2(\mathcal{G}; L^2(\mathbb{P}_X))$, which is isometrically isomorphic to $\mathcal{G} \otimes L^2(\mathbb{P}_X) \cong L^2(\mathbb{P}_X; \mathcal{G})$ (see Remark 5.3). In this case, $f_g = \mathfrak{m}(g)$ by (5.3).

If $\mathfrak{m} \in L_2(\mathcal{G}; \mathcal{H})$, then clearly $f_g = \mathfrak{m}(g) \in \mathcal{H}$ and this shows that (A) \implies (A_{old}).

Now let $[m] \in (\mathcal{G} \otimes \mathcal{H})_{\mathcal{C}}$. Then there exist $h \in \mathcal{G} \otimes \mathcal{H}$ and $c \in \mathcal{G}$ such that $h(x) + c = m(x)$ for \mathbb{P}_X -a.e. $x \in \mathcal{X}$, which implies $f_g = m(g) = h(g) + \langle c, g \rangle_{\mathcal{G}}$ \mathbb{P}_X -a.e. in \mathcal{X} . Since $h(g) \in \mathcal{H}$ and $\langle c, g \rangle_{\mathcal{G}} \in \mathbb{R}$ for each $g \in \mathcal{G}$, this shows that **(B)** \implies **(B_{old})**.

Let $h \in \mathcal{G} \otimes \mathcal{H}$ be such that $[h] = P_{(\mathcal{G} \otimes \mathcal{H})_{\mathcal{C}} L^2_{\mathcal{C}}(\mathbb{P}_X; \mathcal{G})}[m]$. Letting $c := \mathbb{E}[(m - h)(X)] \in \mathcal{G}$ and denoting the unit constant function by $\mathbb{1} \in L^2(\mathbb{P}_X)$, it follows that, for each $h \in \mathcal{H}$ and $g \in \mathcal{G}$,

$$\begin{aligned} 0 &= \langle [h] - [m], [g \otimes h] \rangle_{L^2_{\mathcal{C}}(\mathbb{P}_X; \mathcal{G})} \\ &= \langle h + c \otimes \mathbb{1} - m, g \otimes h - g \otimes \mathbb{E}[h(X)] \rangle_{\mathcal{G} \otimes L^2(\mathbb{P}_X)} \\ &= \langle h(g) + \langle c, g \rangle_{\mathcal{G}} \mathbb{1} - m(g), h - \mathbb{E}[h(X)] \rangle_{L^2(\mathbb{P}_X)} \\ &= \langle [h(g)] - [m(g)], [h] \rangle_{L^2_{\mathcal{C}}(\mathbb{P}_X)}, \end{aligned}$$

where we used (5.3) and $\langle c, g \rangle_{\mathcal{G}} = \mathbb{E}[(m(g) - h(g))(X)]$. Since $h(g) \in \mathcal{H}$, this shows that **(C)** \implies **(C_{old})**.

If $h := P_{\mathcal{G} \otimes \mathcal{H}}^{L^2(\mathbb{P}_X; \mathcal{G})} m \in \mathcal{G} \otimes \mathcal{H}$, then, for each $h \in \mathcal{H}$ and $g \in \mathcal{G}$,

$$0 = \langle h - m, g \otimes h \rangle_{\mathcal{G} \otimes L^2(\mathbb{P}_X)} = \langle h(g) - m(g), h \rangle_{L^2(\mathbb{P}_X)},$$

where we used (5.3). Since $h(g) \in \mathcal{H}$, this shows that **(^uC)** \implies **(^uC_{old})**.

If $m \in \overline{\mathcal{G} \otimes \mathcal{H}}^{L^2(\mathbb{P}_X; \mathcal{G})}$, then there exists a sequence $(h_n)_{n \in \mathbb{N}}$ in $\mathcal{G} \otimes \mathcal{H}$ such that $\|h_n - m\|_{L_2(\mathcal{G}; L^2(\mathbb{P}_X))} \rightarrow 0$ as $n \rightarrow \infty$. Let $g \in \mathcal{G}$ and $h_n := h_n(g) \in \mathcal{H}$, $n \in \mathbb{N}$. Then

$$\|h_n - f_g\|_{L^2(\mathbb{P}_X)} = \|h_n(g) - m(g)\|_{L^2(\mathbb{P}_X)} \leq \|h_n - m\|_{L_2(\mathcal{G}; L^2(\mathbb{P}_X))} \|g\|_{\mathcal{G}} \xrightarrow{n \rightarrow \infty} 0,$$

which proves **(A*)** \implies **(A_{old}*)**.

Finally, let $[m] \in \overline{(\mathcal{G} \otimes \mathcal{H})_{\mathcal{C}}}^{L^2_{\mathcal{C}}(\mathbb{P}_X; \mathcal{G})}$. Then there exists a sequence $(h_n)_{n \in \mathbb{N}}$ in $\mathcal{G} \otimes \mathcal{H}$ such that $\|[h_n] - [m]\|_{\mathcal{G} \otimes L^2_{\mathcal{C}}(\mathbb{P}_X)} \rightarrow 0$ as $n \rightarrow \infty$, that is, $\|h_n + c_n \otimes \mathbb{1} - m\|_{\mathcal{G} \otimes L^2(\mathbb{P}_X)} \rightarrow 0$ as $n \rightarrow \infty$ where $c_n := \mathbb{E}[(m - h_n)(X)] \in \mathcal{G}$ and $\mathbb{1} \in L^2(\mathbb{P}_X)$ is the unit constant function. Let $g \in \mathcal{G}$, $h_n := h_n(g) \in \mathcal{H}$ and $r_n := \langle c_n, g \rangle_{\mathcal{G}} \mathbb{1}$, $n \in \mathbb{N}$. Then, as $n \rightarrow \infty$,

$$\|h_n + r_n - f_g\|_{L^2(\mathbb{P}_X)} = \|h_n(g) + \langle c_n, g \rangle_{\mathcal{G}} \mathbb{1} - m(g)\|_{L^2(\mathbb{P}_X)} \leq \|h_n + c_n \otimes \mathbb{1} - m\|_{L_2(\mathcal{G}; L^2(\mathbb{P}_X))} \|g\|_{\mathcal{G}} \rightarrow 0,$$

which proves **(B*)** \implies **(B_{old}*)**.

The last two implications follow from the above and [22], Theorems 4.1 and 5.1. \square

Proof of Lemma 5.7. Suppose that $(\mathcal{G} \otimes \mathcal{H})_{\mathcal{C}}$ is not dense in $L^2_{\mathcal{C}}(\mathbb{P}_X; \mathcal{G})$. Then there exists $f \in L^2(\mathbb{P}_X; \mathcal{G})$ that is not \mathbb{P}_X -a.e. constant (i.e., there is no $c \in \mathcal{G}$ such that $f(x) = c$ for \mathbb{P}_X -a.e. $x \in \mathcal{X}$) such that $[f] \perp_{L^2_{\mathcal{C}}(\mathbb{P}_X; \mathcal{G})} (\mathcal{G} \otimes \mathcal{H})_{\mathcal{C}}$. Let $\tilde{f} := f - \mathbb{E}[f(X)]$ and $\mathbb{P}_{\tilde{f}} = \tilde{f}_{\#} \mathbb{P}_X$ denote the pushforward measure of \mathbb{P}_X under \tilde{f} .

Then there exists $g_* \in \text{supp}(\mathbb{P}_{\tilde{f}}) \subseteq \mathcal{G}$ such that $g_* \neq 0$ (otherwise $\mathbb{P}_{\tilde{f}}(\mathcal{G} \setminus \{0\}) = 0$ and therefore $\tilde{f} = 0$ \mathbb{P}_X -a.e.). Hence, after proper normalisation of \tilde{f} , we can define the following two *distinct* probability measures on \mathcal{X} :

$$Q_1(E) := \int_E |\langle \tilde{f}(x), g_* \rangle_{\mathcal{G}}| d\mathbb{P}_X(x), \quad Q_2(E) := \int_E |\langle \tilde{f}(x), g_* \rangle_{\mathcal{G}}| - \langle \tilde{f}(x), g_* \rangle_{\mathcal{G}} d\mathbb{P}_X(x)$$

for every measurable subset $E \subseteq \mathcal{X}$. Indeed, for $\varepsilon := \|g_*\|_{\mathcal{G}}/2 > 0$ and any $g = g_* + w \in B_\varepsilon(g_*) := \{g_* + w \mid w \in \mathcal{G}, \|w\|_{\mathcal{G}} < \varepsilon\}$, the reverse triangle inequality implies $\langle g, g_* \rangle_{\mathcal{G}} \geq \|g_*\|_{\mathcal{G}}^2 - \|w\|_{\mathcal{G}} \|g_*\|_{\mathcal{G}} > 2\varepsilon^2$. Hence, since $g_* \in \text{supp}(\mathbb{P}_{\tilde{f}})$, it follows that, for $E = \tilde{f}^{-1}(B_\varepsilon(g_*))$,

$$Q_1(E) - Q_2(E) = \int_E \langle \tilde{f}(x), g_* \rangle_{\mathcal{G}} d\mathbb{P}_X(x) \geq 2\varepsilon^2 \mathbb{P}_X(E) = 2\varepsilon^2 \mathbb{P}_{\tilde{f}}(B_\varepsilon(g_*)) > 0.$$

Since, for every $\mathfrak{h} \in \mathcal{G} \otimes \mathcal{H}$,

$$\langle \tilde{f}, \mathfrak{h} \rangle_{L^2(\mathbb{P}_X; \mathcal{G})} = \langle \tilde{f} - \mathbb{E}[\tilde{f}(X)], \mathfrak{h} \rangle_{L^2(\mathbb{P}_X; \mathcal{G})} \stackrel{[\tilde{f}] \perp (\mathcal{G} \otimes \mathcal{H})^c}{=} \langle \tilde{f} - \mathbb{E}[\tilde{f}(X)], \mathbb{E}[\mathfrak{h}(X)] \rangle_{L^2(\mathbb{P}_X; \mathcal{G})} = 0,$$

it follows that $\tilde{f} \perp_{L^2(\mathbb{P}_X; \mathcal{G})} \mathcal{G} \otimes \mathcal{H}$. Let $Z_1 \sim Q_1$, $Z_2 \sim Q_2$ and $x \in \mathcal{X}$. Since $g_* \otimes \varphi(x) \in \mathcal{G} \otimes \mathcal{H}$,

$$(\mathbb{E}[\varphi(Z_1)] - \mathbb{E}[\varphi(Z_2)])(x) = \int_{\mathcal{X}} k(x, x') \langle \tilde{f}(x'), g_* \rangle_{\mathcal{G}} d\mathbb{P}_X(x') = \langle g_* \otimes \varphi(x), \tilde{f} \rangle_{L^2(\mathbb{P}_X; \mathcal{G})} = 0,$$

where we used (5.2), contradicting the assumption of k being characteristic. \square

Proof of Theorem 5.8. By Assumption 5.4(B), $\mathfrak{m}(x) = \tilde{f}(x) + c$ with $\tilde{f} \in \mathcal{G} \otimes \mathcal{H}$ and $c \in \mathcal{G}$. As discussed in Remark 5.3, we can view $\tilde{f} \in \mathcal{G} \otimes \mathcal{H}$ as an element of both $L^2(\mathbb{P}_X; \mathcal{G})$ and $L_2(\mathcal{H}; \mathcal{G})$, and thereby \mathfrak{m} as an element of $A_2(\mathcal{H}; \mathcal{G})$. Hence, (5.2) and the injectivity of φ imply that

$$\mathbb{E}[U|V] = \mathbb{E}[\psi(Y)|X] = \mathfrak{m}(X) = \tilde{f}(X) + c = \tilde{f}(\varphi(X)) + c = \mathfrak{m} \circ V \quad \text{a.s.}$$

Since $\mathfrak{m} \in A_2(\mathcal{H}; \mathcal{G}) \subseteq A(\mathcal{H}; \mathcal{G})$, the statements follow from Theorem 4.8; the inclusion $\text{ran } C_{VU} \subseteq \text{ran } C_V$ follows from Assumption 5.4(B) by Proposition 5.6, cf. Figure 3. \square

Proof of Theorem 5.11. By Assumption 5.4(B*), there exists a sequence $\tilde{f}^{(n)} \in \mathcal{G} \otimes \mathcal{H}$, $n \in \mathbb{N}$, such that

$$\|\tilde{f}^{(n)} - \mathfrak{m}\|_{L^2_{\mathcal{C}}(\mathbb{P}_X; \mathcal{G})} \xrightarrow{n \rightarrow \infty} 0.$$

Therefore, denoting $c^{(n)} := \mathbb{E}[\mathfrak{m}(X) - \tilde{f}^{(n)}(X)] \in \mathcal{G}$,

$$\|\tilde{f}^{(n)}(X) + c^{(n)} - \mathfrak{m}(X)\|_{L^2(\mathbb{P}; \mathcal{G})} \xrightarrow{n \rightarrow \infty} 0.$$

As discussed in Remark 5.3, $\tilde{f}^{(n)} \in \mathcal{G} \otimes \mathcal{H}$ can be seen as an element of both $L^2(\mathbb{P}_X; \mathcal{G})$ and $L_2(\mathcal{H}; \mathcal{G})$. Hence, (5.2) and the injectivity of φ imply

$$\begin{aligned} \mathbb{E}[U|V] &= \mathbb{E}[\psi(Y)|X] = \mathfrak{m}(X) = \lim_{n \rightarrow \infty} \tilde{f}^{(n)}(X) + c^{(n)} \\ &= \lim_{n \rightarrow \infty} \tilde{f}^{(n)}(\varphi(X)) + c^{(n)} = \lim_{n \rightarrow \infty} \gamma^{(n)} \circ V \quad \text{a.s.,} \end{aligned}$$

where the limits are in $L^2(\mathbb{P}; \mathcal{G})$ and $\gamma^{(n)}(h) := \tilde{f}^{(n)}(h) + c^{(n)}$. Since $\gamma^{(n)} \in A_2(\mathcal{H}; \mathcal{G}) \subseteq A(\mathcal{H}; \mathcal{G})$, this implies that $\mathbb{E}[U|V] \in \overline{A_{\mathcal{G}} \circ V}$ and thereby $\mathbb{E}[U|V] = \mathbb{E}^A[U|V]$. The claim now follows from Theorem 4.13. \square

Appendix: Technical results

The following well-known result due to [9], Theorem 1, (see also [14], Theorem 2.1) is used several times.

Theorem A.1. *Let \mathcal{H} , \mathcal{H}_1 , and \mathcal{H}_2 be Hilbert spaces and let $A: \mathcal{H}_1 \rightarrow \mathcal{H}$ and $B: \mathcal{H}_2 \rightarrow \mathcal{H}$ be bounded linear operators with $\text{ran } A \subseteq \text{ran } B$. Then $Q := B^\dagger A: \mathcal{H}_1 \rightarrow \mathcal{H}_2$ is a well-defined and bounded linear operator, where B^\dagger denotes the Moore–Penrose pseudo-inverse of B . It is the unique operator that satisfies the conditions*

$$A = BQ, \quad \ker Q = \ker A, \quad \text{ran } Q \subseteq \overline{\text{ran } B^*}. \quad (\text{A.1})$$

Remark A.2. In the original work of [9] only the existence of a bounded operator Q such that $A = BQ$ was shown. However, the construction of Q in the proof is identical to that of B^\dagger (multiplied by A). This connection has been observed before by [1], Corollary 2.2 and Remark 2.3, where it was proven in the case of closed range operators, leaving the proof of the general case to the reader. Moreover, [9], Theorem 1, only treats the case $\mathcal{H} = \mathcal{H}_1 = \mathcal{H}_2$; the general case is mentioned as a remark at the end of his paper.

Further, we are going to use the following characterisations of cross-covariance operators due to [3], Theorem 1.

Theorem A.3. *Under the notation of Section 3, there exists a unique bounded linear operator $R_{VU}: \mathcal{G} \rightarrow \mathcal{H}$ with operator norm $\|R_{VU}\| \leq 1$ such that*

$$C_{VU} = C_V^{1/2} R_{VU} C_U^{1/2}, \quad R_{VU} = P_{(\ker C_V)^\perp} R_{VU} P_{(\ker C_U)^\perp}. \quad (\text{A.2})$$

Remark A.4. If $\mathcal{H} = \mathcal{G} = \mathbb{R}$, then R_{VU} coincides with the Pearson correlation coefficient.

This paper makes extensive use of the following two basic results.

Lemma A.5. *Let $A: \mathcal{H} \rightarrow \mathcal{G}$ be a trace-class operator such that $\text{tr}(AB^*) = 0$ for any bounded operator $B \in \mathcal{L}(\mathcal{H}; \mathcal{G})$. Then $A = 0$.*

Proof. Choosing $B = A$ yields $\|A\|_{\mathcal{L}_2} = \text{tr}(AA^*)^{1/2} = 0$, hence $A = 0$. □

Lemma A.6. *With the notation of Section 3, let $U' \in L^2(\Omega, \Sigma, \mathbb{P}; \mathcal{G})$ and $\gamma \in \mathbf{A}_V(\mathcal{H}; \mathcal{G})$. Then*

- (a) $\langle U - \mu_U, U' \rangle_{L^2(\mathbb{P}; \mathcal{G})} = \text{tr}(\text{Cov}[U, U'])$;
- (b) $\text{Cov}[\gamma \circ V, U] = \overline{\gamma} C_{VU}$ and $\text{Cov}[U, \gamma \circ V] = (\overline{\gamma} C_{VU})^*$.

If $\gamma \in \mathbf{A}(\mathcal{H}; \mathcal{G})$, then the last equation can be simplified to $\text{Cov}[U, \gamma \circ V] = C_{UV} \overline{\gamma}^$.*

Proof. Let $(e_j)_{j \in \mathcal{J}}$, $\mathcal{J} \subseteq \mathbb{N}$ be an orthonormal basis of \mathcal{G} . Then

$$\begin{aligned} \langle U - \mu_U, U' \rangle_{L^2(\mathbb{P}; \mathcal{G})} &= \langle U - \mu_U, U' - \mu_{U'} \rangle_{L^2(\mathbb{P}; \mathcal{G})} \\ &= \mathbb{E}[\langle U - \mu_U, U' - \mu_{U'} \rangle_{\mathcal{G}}] \\ &= \sum_{j \in \mathcal{J}} \mathbb{E}[\langle e_j, U - \mu_U \rangle_{\mathcal{G}} \langle U' - \mu_{U'}, e_j \rangle_{\mathcal{G}}] \end{aligned}$$

$$\begin{aligned}
&= \sum_{j \in \mathcal{J}} \langle e_j, C_{UU'} e_j \rangle_{\mathcal{G}} \\
&= \text{tr}[\text{Cov}[U, U']],
\end{aligned}$$

proving (a). Since $V \in L^2(\mathbb{P}; \mathcal{H})$ and $\gamma \circ V \in L^2(\mathbb{P}; \mathcal{G})$, all covariance operators are well defined and so, for $g \in \mathcal{G}$,

$$\text{Cov}[\gamma V, U](g) = \mathbb{E}[\gamma(V - \mu_V)\langle U - \mu_U, g \rangle_{\mathcal{G}}] = \gamma \mathbb{E}[(V - \mu_V)\langle U - \mu_U, g \rangle_{\mathcal{G}}] = \gamma \text{Cov}[V, U](g),$$

proving (b). □

Funding

IK and TJS are supported in part by the Deutsche Forschungsgemeinschaft (DFG) through project TrU-2 “Demand modelling and control for e-commerce using RKHS transfer operator approaches” of the Excellence Cluster “MATH+ The Berlin Mathematics Research Centre” (EXC-2046/1, project 390685689). TJS is further supported by the DFG project 415980428. BS has been supported by the DFG project 389483880.

References

- [1] Arias, M.L., Corach, G. and Gonzalez, M.C. (2008). Generalized inverses and Douglas equations. *Proc. Amer. Math. Soc.* **136** 3177–3183. [MR2407082](#) <https://doi.org/10.1090/S0002-9939-08-09298-8>
- [2] Aubin, J.-P. (2000). *Applied Functional Analysis*, 2nd ed. *Pure and Applied Mathematics (New York)*. New York: Wiley Interscience. [MR1782330](#) <https://doi.org/10.1002/9781118032725>
- [3] Baker, C.R. (1973). Joint measures and cross-covariance operators. *Trans. Amer. Math. Soc.* **186** 273–289. [MR0336795](#) <https://doi.org/10.2307/1996566>
- [4] Brémaud, P. (2017). *Discrete Probability Models and Methods: Probability on Graphs and Trees, Markov Chains and Random Fields, Entropy and Coding. Probability Theory and Stochastic Modelling* **78**. Cham: Springer. [MR3616988](#) <https://doi.org/10.1007/978-3-319-43476-6>
- [5] Chatterji, S.D. (1960). Martingales of Banach-valued random variables. *Bull. Amer. Math. Soc.* **66** 395–398. [MR0119242](#) <https://doi.org/10.1090/S0002-9904-1960-10471-5>
- [6] Chilès, J.-P. and Delfiner, P. (2012). *Geostatistics: Modeling Spatial Uncertainty*, 2nd ed. *Wiley Series in Probability and Statistics*. Hoboken, NJ: Wiley. [MR2850475](#) <https://doi.org/10.1002/9781118136188>
- [7] Corach, G., Maestripieri, A. and Stojanoff, D. (2001). Oblique projections and Schur complements. *Acta Sci. Math. (Szeged)* **67** 337–356. [MR1830147](#)
- [8] Diestel, J. and Uhl, J.J. Jr. (1977). *Vector Measures*. Providence, RI: Amer. Math. Soc. [MR0453964](#)
- [9] Douglas, R.G. (1966). On majorization, factorization, and range inclusion of operators on Hilbert space. *Proc. Amer. Math. Soc.* **17** 413–415. [MR0203464](#) <https://doi.org/10.2307/2035178>
- [10] Dudley, R.M. (2002). *Real Analysis and Probability. Cambridge Studies in Advanced Mathematics* **74**. Cambridge: Cambridge Univ. Press. [MR1932358](#) <https://doi.org/10.1017/CBO9780511755347>
- [11] Engl, H.W., Hanke, M. and Neubauer, A. (1996). *Regularization of Inverse Problems. Mathematics and Its Applications* **375**. Dordrecht: Kluwer Academic. [MR1408680](#)
- [12] Ernst, O.G., Sprungk, B. and Starkloff, H.-J. (2015). Analysis of the ensemble and polynomial chaos Kalman filters in Bayesian inverse problems. *SIAM/ASA J. Uncertain. Quantificat.* **3** 823–851. [MR3400030](#) <https://doi.org/10.1137/140981319>
- [13] Evensen, G. (2009). *Data Assimilation: The Ensemble Kalman Filter*, 2nd ed. Berlin: Springer. [MR2555209](#) <https://doi.org/10.1007/978-3-642-03711-5>

- [14] Fillmore, P.A. and Williams, J.P. (1971). On operator ranges. *Adv. Math.* **7** 254–281. [MR0293441](#) [https://doi.org/10.1016/S0001-8708\(71\)80006-3](https://doi.org/10.1016/S0001-8708(71)80006-3)
- [15] Folland, G.B. (1999). *Real Analysis: Modern Techniques and Their Applications*, 2nd ed. *Pure and Applied Mathematics (New York)*. New York: Wiley. [MR1681462](#)
- [16] Fukumizu, K., Song, L. and Gretton, A. (2013). Kernel Bayes' rule: Bayesian inference with positive definite kernels. *J. Mach. Learn. Res.* **14** 3753–3783. [MR3159407](#)
- [17] Goldstein, M. (1999). Bayes linear analysis. In *Encyclopaedia of Statistical Sciences* (S. Kotz, B.C. Read, N. Balakrishnan, B. Vidakovic and N.L. Johnson, eds.) 29–34. Chichester: Wiley.
- [18] Goldstein, M. and Wooff, D. (2007). *Bayes Linear Statistics: Theory and Methods*. *Wiley Series in Probability and Statistics*. Chichester: Wiley. [MR2335584](#) <https://doi.org/10.1002/9780470065662>
- [19] Hairer, M., Stuart, A.M., Voss, J. and Wiberg, P. (2005). Analysis of SPDEs arising in path sampling. I. The Gaussian case. *Commun. Math. Sci.* **3** 587–603. [MR2188686](#)
- [20] Hastie, T., Tibshirani, R. and Friedman, J. (2009). *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*, 2nd ed. *Springer Series in Statistics*. New York: Springer. [MR2722294](#) <https://doi.org/10.1007/978-0-387-84858-7>
- [21] Kallenberg, O. (2006). *Foundations of Modern Probability. Probability and Its Applications (New York)*. New York: Springer. [MR1464694](#)
- [22] Klebanov, I., Schuster, I. and Sullivan, T.J. (2020). A rigorous theory of conditional mean embeddings. *SIAM J. Math. Data Sci.* **2** 583–606. [MR4121886](#) <https://doi.org/10.1137/19M1305069>
- [23] Klenke, A. (2013). *Wahrscheinlichkeitstheorie*, 3rd ed. Berlin: Springer. <https://doi.org/10.1007/978-3-642-36018-3>
- [24] Klus, S., Husic, B.E., Mollenhauer, M. and Noé, F. (2019). Kernel methods for detecting coherent structures in dynamical data. *Chaos* **29** 123112, 15. [MR4041090](#) <https://doi.org/10.1063/1.5100267>
- [25] Klus, S., Nüske, F., Koltai, P., Wu, H., Kevrekidis, I., Schütte, C. and Noé, F. (2018). Data-driven model reduction and transfer operator approximation. *J. Nonlinear Sci.* **28** 985–1010. [MR3800253](#) <https://doi.org/10.1007/s00332-017-9437-7>
- [26] Klus, S., Schuster, I. and Muandet, K. (2020). Eigendecompositions of transfer operators in reproducing kernel Hilbert spaces. *J. Nonlinear Sci.* **30** 283–315. [MR4054854](#) <https://doi.org/10.1007/s00332-019-09574-z>
- [27] Mandelbaum, A. (1984). Linear estimators and measurable linear transformations on a Hilbert space. *Z. Wahrsch. Verw. Gebiete* **65** 385–397. [MR0731228](#) <https://doi.org/10.1007/BF00533743>
- [28] Meise, R. and Vogt, D. (1997). *Introduction to Functional Analysis. Oxford Graduate Texts in Mathematics 2*. New York: The Clarendon Press, Oxford Univ. Press. [MR1483073](#)
- [29] Owahdi, H. and Scovel, C. (2018). Conditioning Gaussian measure on Hilbert space. *J. Math. Stat. Anal.* **1** 1–15.
- [30] Sazonov, V. (1958). On characteristic functionals. *Teor. Veroyatn. Primen.* **3** 201–205. [MR0098423](#)
- [31] Schillings, C. and Stuart, A.M. (2017). Analysis of the ensemble Kalman filter for inverse problems. *SIAM J. Numer. Anal.* **55** 1264–1290. [MR3654885](#) <https://doi.org/10.1137/16M105959X>
- [32] Schwantes, C.R. and Pande, V.S. (2015). Modeling molecular kinetics with tICA and the kernel trick. *J. Chem. Theory Comput.* **11** 600–608. <https://doi.org/10.1021/ct5007357>
- [33] Song, L., Huang, J., Smola, A. and Fukumizu, K. (2009). Hilbert space embeddings of conditional distributions with applications to dynamical systems. In *Proceedings of the 26th Annual International Conference on Machine Learning* 961–968. <https://doi.org/10.1145/1553374.1553497>
- [34] Stein, M.L. (1999). *Interpolation of Spatial Data: Some Theory for Kriging. Springer Series in Statistics*. New York: Springer. [MR1697409](#) <https://doi.org/10.1007/978-1-4612-1494-6>
- [35] Steinwart, I. and Christmann, A. (2008). *Support Vector Machines. Information Science and Statistics*. New York: Springer. [MR2450103](#)
- [36] Tarieladze, V. and Vakhania, N. (2007). Disintegration of Gaussian measures and average-case optimal algorithms. *J. Complexity* **23** 851–866. [MR2371996](#) <https://doi.org/10.1016/j.jco.2007.04.005>